

May 1985

Report No. STAN-CS-85-1051



PB96-146998

Special Relations in Automated Deduction

by

Zohar Manna

Richard Waldinger

DISTRIBUTION STATEMENT A

Approved for public release
Distribution Unlimited

Department of Computer Science

Stanford University
Stanford, CA 94305

19970609 037



PRINT QUALITY INSPECTED 8.

SPECIAL RELATIONS IN AUTOMATED DEDUCTION

Zohar Manna
Computer Science Department
Stanford University

Richard Waldinger
Artificial Intelligence Center
SRI International

ABSTRACT

Two deduction rules are introduced to give streamlined treatment to relations of special importance in an automated theorem-proving system. These rules, the *relation replacement* and *relation matching* rules, generalize to an arbitrary binary relation the paramodulation and E-resolution rules, respectively, for equality, and may operate within a nonclausal or clausal system. The new rules depend on an extension of the notion of *polarity* to apply to subterms as well as to subsentences, with respect to a given binary relation. The rules allow us to eliminate troublesome axioms, such as transitivity and monotonicity, from the system; proofs are shorter and more comprehensible, and the search space is correspondingly deflated.

1. INTRODUCTION

In any theorem-proving system, the task of representing properties of objects is shared between axioms and rules of inference. The axioms of the system are easier to introduce and modify, because they are expressed in a logical language. However, because axioms are declarative rather than imperative, they are given no individual heuristic controls. The rules of inference, on the other hand, cannot be altered without reprogramming the system, and they are usually expressed in the system's programming language. However, the rules can be given individual heuristic controls and strategies.

It is customary to use rules of inference to express properties of the logical connectives, which are the same from one theory to the next, and to use axioms to express properties of constants, functions, and relations, which may vary. It is hazardous, however, to express certain properties of functions and relations by axioms. Some properties of the equality relation, for example, are rarely represented axiomatically. For one thing, in a first-order system indefinitely many axioms are necessary to represent the substitutivity property of this relation, depending on how many function and relation symbols are in the vocabulary of the theory.

For instance, for a binary function symbol $f(x, y)$, we must introduce two *functional-substitutivity* axioms,

$$\begin{array}{ll} \text{if } x = y & \text{if } x = y \\ \text{then } f(x, z) = f(y, z) & \text{and } \text{then } f(z, x) = f(z, y), \end{array}$$

and for a binary predicate symbol $p(x, y)$, we must introduce two *predicate-substitutivity* axioms,

$$\begin{array}{ll} \text{if } x = y & \text{if } x = y \\ \text{then if } p(x, z) \text{ then } p(y, z) & \text{and } \text{then if } p(z, x) \text{ then } p(z, y). \end{array}$$

An abbreviated version of this paper appears in the proceedings of the Twelfth International Colloquium on Automata, Languages, and Programming (ICALP), Nafplion, Greece, July 1985.

This research was supported in part by the National Science Foundation under grants MCS-82-14523 and MCS-81-05565, by the Defense Advanced Research Projects Agency under contract N00039-84-C-0211, by the United States Air Force Office of Scientific Research under contract AFOSR-81-0014, by the Office of Naval Research under contract N00014-84-C-0706, and by a contract from the International Business Machines Corporation.

(We tacitly quantify variables universally over the entire sentence.) In general, for each n -ary function symbol $f(x_1, \dots, x_n)$, we introduce n *functional-substitutivity* axioms. Similarly, for each n -ary predicate symbol $p(x_1, \dots, x_n)$, we introduce n *predicate-substitutivity* axioms.

More importantly, axioms for equality are difficult to control strategically, because they have many irrelevant consequences. An axiom such as *transitivity*,

$$\begin{array}{l} \text{if } x = y \text{ and } y = z \\ \text{then } x = z, \end{array}$$

will allow us to derive logical consequences from any sentence mentioning the equality relation. Few of these consequences will have any bearing on the proof.

In response to this problem, some theorem-proving researchers have paraphrased their theories to avoid explicit mention of the equality axiom (e.g., Kowalski [79]). Others have adopted special inference rules for dealing with equality. In resolution systems, two equality rules, paramodulation (Wos and Robinson [69]) and E-resolution (Morris [69]) have been found to be effective. Variations of these rules are used in many theorem provers today (e.g., Boyer and Moore [79], Digricoli [83]). By a single application of either of these rules, we can derive conclusions that would require several steps if the properties of equality were represented axiomatically. The proofs are markedly shorter, and the search spaces are even more dramatically compressed because the axioms and intermediate steps are not required. Within their limited domain of application, theorem-proving systems using these rules surpass most human beings in their capabilities.

SPECIAL RELATIONS

The authors became involved in theorem proving because of its application to program synthesis, the derivation of a program to meet a given specification. We have been pursuing a deductive approach to this problem, under which computer programming is regarded as a theorem-proving task. In the proofs required for program synthesis, certain relations assume special importance. Again and again, proofs require us to reason not only about the equality relation, but also about the less-than relation $<$ (over the integers or reals), the subset relation \subseteq , the sublist relation \leq_{list} , or the subtree relation \leq_{tree} . To represent the transitivity and other properties of these relations axiomatically leads to many of the same problems that were faced in dealing with equality: the axioms apply almost everywhere, spawning innumerable consequences that swamp the system. Yet we would not want to implement a new inference rule for each of the relations we find important.

Both the paramodulation and the E-resolution rules are based on the *substitutivity* property of equality, that if two elements are equal they may be used interchangeably; i.e., for any sentence $P(x, y)$, the sentence

$$\begin{array}{l} \text{if } x = y \\ \text{then if } P(x, y) \text{ then } P(y, x) \end{array}$$

is valid. Here $P(y, x)$ is the result of replacing in $P(x, y)$ certain (perhaps none) of the occurrences of x with y , and certain (perhaps none) of the occurrences of y with x . (The notations we use here informally will be defined systematically later on. We assume throughout that sentences are quantifier-free.)

We observe that many of the relations we regard as important exhibit substitutivity properties similar to the above property of equality, but under restricted circumstances. For example, over the nonnegative integers, we can show that

$$\begin{array}{l} \text{if } x < y \\ \text{then if } a \leq x \cdot b \\ \text{then } a \leq y \cdot b \end{array}$$

and, over the lists, we can show that

$$\begin{array}{l} \text{if } x \leq_{list} y \\ \text{then if } u \in x \\ \quad \text{then } u \in y. \end{array}$$

Knowing that $x < y$ or that $x \leq_{list} y$ does not allow us to use x and y interchangeably, but it does allow us to replace certain occurrences of x with y , and vice versa.

Based on such substitutivity properties, we can introduce two deduction rules that generalize the paramodulation and E-resolution rules for equality to an arbitrary relation, under appropriate circumstances. Just as the equality rules enable us to drop the transitivity and substitutivity axioms for equality, the new relation rules enable us to drop the corresponding troublesome axioms for the relations of our theory.

POLARITY

For the equality relation, knowing that $x = y$ allows us to replace in a given sentence any occurrence of x with y and any occurrence of y with x , obtaining a sentence that follows from the given one. For an arbitrary binary relation \prec , knowing that $x \prec y$ still may allow us to replace certain occurrences of x with y and certain occurrences of y with x . We describe a syntactic procedure that, for a given relation \prec , identifies which occurrences of x and y in a given sentence can be replaced, provided we know that $x \prec y$.

More precisely, we identify particular occurrences of subexpressions of a given sentence as being positive (+), negative (-), or both, or neither, with respect to \prec . If $x \prec y$, positive occurrences of x can be replaced with y , and negative occurrences of y can be replaced with x . In other words, we can establish the substitutivity property that, for any sentence $P(x^+, y^-)$, the sentence

$$\begin{array}{l} \text{if } x \prec y \\ \text{then if } P(x^+, y^-) \text{ then } P(y^+, x^-) \end{array}$$

is valid (over the theory in question). Here $P(y^+, x^-)$ is the sentence obtained from $P(x^+, y^-)$ by replacing certain positive occurrences of x with y and certain negative occurrences of y with x . With respect to the equality relation, every subexpression is both positive and negative; therefore, if we take \prec to be $=$, this property reduces to the substitutivity of equality.

Our new rules are based on the above substitutivity property just as the equality rules are based on the substitutivity of equality. The new rules, like the equality rules, allow us to perform in a single application inferences that would require many steps in a conventional system. Proofs are shorter and closer to an intuitive argument; the search space is condensed accordingly.

NONCLAUSAL DEDUCTION

The paramodulation and E-resolution rules are formulated for sentences in clausal form (a disjunction of atomic sentences and their negations); on the other hand, the two corresponding rules we introduce apply to free-form sentences, with a full set of logical connectives (cf. Manna and Waldinger [80], Murray [82], Stickel [82]). By adopting such a nonclausal system, we avoid the proliferation of sentences and the disintegration of intuition that accompany the translation to clausal form. Also, it is awkward to express the mathematical induction principle in a clausal system, because we must do induction on sentences that may require more than one clause to express. On the other hand, our rules are also immediately and directly applicable to clausal theorem-proving systems.

OUTLINE

In the following section, **Preliminaries**, we sketch the basic concepts of logic that we use in this paper and we briefly outline a nonclausal deduction system. Readers who are familiar with this material should skim the section anyway, to become acquainted with our terminology and notations.

In **Relational Polarity** we introduce our central notion, the polarity of a subexpression of a sentence with respect to a given relation.

We then describe, in **The Relation Replacement Rule**, a new deduction rule that allows us to replace a subexpression of a sentence with another expression, under a wide variety of circumstances. This is our generalization of the paramodulation rule.

The rules in our system can be applied when two subexpressions can be unified. However, our second deduction rule, described in **The Relation Matching Rule**, allows us to draw a conclusion even though two subexpressions fail to unify. (Typically this rule is applied when the two subexpressions "nearly" unify.) This is our generalization of the E-resolution rule.

In **Strengthening** we tighten up our theory of polarity to allow the relation replacement rule to draw a stronger conclusion, in many circumstances.

In **Extensions**, we indicate how the notions in this paper can be extended to apply to sentences which contain explicit quantifiers and to define polarity with respect to functions as well as relations; we develop more general, conditional versions of all the rules; and we show how our results apply to problems in automated planning.

2. PRELIMINARIES

Before we can define our central notion, that of polarity of a subexpression with respect to a relation, we must introduce some concepts and notations. We will be brief and informal, because we believe that this material will be familiar to most readers.

EXPRESSIONS

We consider *terms* composed (in the usual way) of the following symbols:

- The constant symbols $a, b, c, a_1, \dots, s, t$, and special constants such as 0.
- The variable symbols $u, v, w, x, y, u_1, \dots$
- The n -ary function symbols f, g, h, f_1, \dots and special symbols such as $+$.

Thus $a, x, f(a, x)$, and $f(a, x) + 0$ are terms.

We consider *propositions* composed (in the usual way) from terms and the following symbols:

- The truth symbols (logical constants) *true* and *false*.
- The n -ary relation symbols p, q, r, p_1, \dots and special symbols such as $=$ and $<$.

Thus *true* and $p(a, g(x))$ are propositions.

We consider *sentences* composed (in the usual way) from propositions and the following symbols:

- The logical connectives *not*, *and*, *or*, *if-then*, \equiv (*if-and-only-if*), *if-then-else*.

Thus $(a < 0)$ or $\text{not}(p(a, g(x)))$ is a sentence.

The *operators* consist of the function and the relation symbols. The *expressions* consist of the terms and the sentences. Note that we do not include the quantifiers \forall and \exists in our language. The *ground* expressions are those that contain no variables. The expressions that occur in a given expression are its *subexpressions*. They are said to be *proper* if they are distinct from the entire expression.

REPLACEMENT

We introduce the operation of replacing subexpressions of a given expression with other expressions. We actually have two distinct notions of replacement, depending on whether or not every occurrence of the subexpression is to be replaced.

Suppose s , t , and e are expressions, where s and t are either both sentences or both terms. If we write e as $e[s]$, then $e[t]$ denotes the expression obtained by replacing every occurrence of s in $e[s]$ with t ; we call this a *total replacement*. If we write e as $e\langle s \rangle$, then $e\langle t \rangle$ denotes the expression obtained by replacing certain (perhaps none) of the occurrences of s in $e\langle s \rangle$ with t ; we call this a *partial replacement*.

When we say we replace certain (perhaps none) of the occurrences of s , we mean that we replace zero, one, or more occurrences. We do not require that $e[s]$ or $e\langle s \rangle$ actually contain any occurrences of s ; if not, $e[t]$ and $e\langle t \rangle$ are the same as $e[s]$ and $e\langle s \rangle$, respectively. Also, while the result of a total replacement is unique, a partial replacement can produce any of several expressions.

For example, if $e[s]$ is $p(s, s, b)$, then $e[t]$ is $p(t, t, b)$. On the other hand, if $e\langle s \rangle$ is $p(s, s, b)$, then $e\langle t \rangle$ could be any of $p(s, s, b)$, $p(t, s, b)$, $p(s, t, b)$, or $p(t, t, b)$. If we want to be more specific about which occurrences are replaced, we must do so in words.

A partial replacement is *invertible*, in the sense that any sentence $e\langle s \rangle$ can be retrieved by replacing certain occurrences of t in $e\langle t \rangle$ with s . The occurrences of t to be replaced are precisely the ones introduced in obtaining $e\langle t \rangle$ in the first place. For example, if $e\langle s \rangle$ is $p(s, s, t)$, and $e\langle t \rangle$ is $p(s, t, t)$, then $e\langle s \rangle$ can be retrieved by replacing the newly introduced occurrence of t in $e\langle t \rangle$ with s .

Total replacement, on the other hand, is not invertible in the same sense. For example, if $e[s]$ is $p(s, s, t)$, then $e[t]$ is $p(t, t, t)$, and $e[s]$ cannot be obtained from $e[t]$ by replacing every occurrence of t in $e[t]$ with s .

MULTIPLE REPLACEMENT

We can extend the definition to allow the replacement of several subexpressions at once:

Suppose $s_1, \dots, s_n, t_1, \dots, t_n$, and e are expressions, where the s_i are distinct and, for each i , s_i and t_i are either both sentences or both terms. If we write e as $e[s_1, \dots, s_n]$, then $e[t_1, \dots, t_n]$ denotes the expression obtained by replacing simultaneously every occurrence of each expression s_i in e with the corresponding expression t_i ; we call this a *multiple total replacement*. If we write e as $e\langle s_1, \dots, s_n \rangle$, then $e\langle t_1, \dots, t_n \rangle$ denotes any of the expressions obtained by replacing simultaneously certain (perhaps none) of the occurrences of some of the expressions s_i in e with the corresponding expression t_i ; we call this a *multiple partial replacement*.

The replacements are made simultaneously in a single stage. For example, if $e[a, b]$ is $f(a, b)$, then $e[b, c]$ is $f(b, c)$. On the other hand, if $e\langle a, b \rangle$ is $f(a, b)$, then $e\langle b, c \rangle$ could denote any of $f(a, b)$, $f(b, b)$, $f(a, c)$, or $f(b, c)$. Even though a is replaced by b and b is replaced by c , the newly introduced occurrences of b are not replaced by c .

The replacements are made from the top down. For example, if $e[p(a, b), a]$ is $p(a, b)$, then $e[\text{true}, b]$ is true . We replace both $p(a, b)$ and a , but a is a subexpression of $p(a, b)$. In such cases, by convention, it is the

outermost subexpression that is replaced. (For the corresponding partial replacement, either subexpression can be replaced.)

By attaching a numerical superscript, we can specify exactly how many subexpression occurrences are to be replaced in a partial replacement. Suppose $s_1, \dots, s_n, t_1, \dots, t_n$, and $e(s_1, \dots, s_n)$ are expressions and k is a nonnegative integer, where the s_i are distinct and, for each i , s_i and t_i are either both sentences or both terms. Then $e(t_1, \dots, t_n)^k$ is the result of replacing in $e(s_1, \dots, s_n)$ precisely k occurrences of s_1, \dots, s_n with the corresponding expression t_1, \dots, t_n . [We assume that at least k occurrences exist.]

Note that precisely k occurrences are replaced altogether. For example, suppose $e(a, b)$ is $c \leq f(a, a, b)$; then $e(a+1, b+1)^2$ could denote any of

$$c \leq f(a+1, a+1, b), \quad c \leq f(a+1, a, b+1), \quad \text{or} \quad c \leq f(a, a+1, b+1),$$

but not

$$c \leq f(a+1, a+1, b+1) \quad \text{or} \quad c \leq f(a+1, a, b).$$

We may also write $e(t_1, t_2, \dots, t_n)^{k, \ell}$ to indicate that precisely k or ℓ replacements are made in the expression $e(s_1, s_2, \dots, s_n)$.

SUBSTITUTIONS

We have a special notation for a substitution, indicating the total replacement of variables with terms. A theory of substitutions was developed by Robinson [65], in the paper in which the resolution principle was introduced. A fuller exposition of this theory appears in Manna and Waldinger [81].

For any distinct variables x_1, x_2, \dots, x_n and any terms t_1, t_2, \dots, t_n , a *substitution*

$$\theta : \{x_1 \leftarrow t_1, x_2 \leftarrow t_2, \dots, x_n \leftarrow t_n\}$$

is a set of replacement pairs $x_i \leftarrow t_i$. Note that there are no substitutions of form $\{x \leftarrow a, x \leftarrow b, \dots\}$, where a and b are distinct. (If a and b are identical, then the set $\{x \leftarrow a, x \leftarrow a, \dots\}$ is the same as the set $\{x \leftarrow a, \dots\}$.) The *empty substitution* $\{\}$ is the set of no replacement pairs.

For any substitution θ and expression e , we denote by $e\theta$ the expression obtained by *applying* θ to e , i.e., by simultaneously replacing every occurrence of the variable x_i in e with the expression t_i , for each replacement pair $x_i \leftarrow t_i$ in θ . We also say that $e\theta$ is an *instance* of e . For example,

$$p(x, y)\{x \leftarrow y, y \leftarrow a\} = p(y, a).$$

The empty substitution $\{\}$ has the property that $e\{\} = e$ for any expression e .

Two substitutions θ and λ are said to be *equal* if they have the same effect on any expression, i.e., if, for any expression e ,

$$e\theta = e\lambda.$$

For example,

$$\{x \leftarrow a, y \leftarrow b\} = \{x \leftarrow a, y \leftarrow b, z \leftarrow z\}.$$

Two substitutions θ and λ are equal if they agree on all variables, i.e., if $x\theta = x\lambda$ for all variables x .

For any variable x , term t , and substitution θ , the result

$$(x \leftarrow t) \circ \theta$$

of adding the replacement pair $x \leftarrow t$ to θ is defined to be the substitution that replaces x with t but agrees with θ on all other variables. It is thus defined by the properties

$$x((x \leftarrow t) \circ \theta) = t$$

$$y((x \leftarrow t) \circ \theta) = y\theta, \text{ for all variables } y \text{ distinct from } x.$$

Note that θ may already replace x with some term t' ; if so, that replacement is superseded by the new one.

For example,

$$(y \leftarrow b) \circ \{ \} = \{y \leftarrow b\}$$

$$(x \leftarrow a) \circ \{y \leftarrow b\} = \{x \leftarrow a, y \leftarrow b\}$$

$$(y \leftarrow c) \circ \{x \leftarrow a, y \leftarrow b\} = \{x \leftarrow a, y \leftarrow c\}$$

$$(x \leftarrow x) \circ \{x \leftarrow a\} = \{ \}.$$

We write $(x \leftarrow t) \circ (y \leftarrow t') \circ \theta$ as an abbreviation for $(x \leftarrow t) \circ ((y \leftarrow t') \circ \theta)$.

The composition $\theta\lambda$ of two substitutions θ and λ is defined by the properties

$$\{ \}\lambda = \lambda$$

$$((x \leftarrow t) \circ \theta)\lambda = (x \leftarrow t\lambda) \circ (\theta\lambda)$$

for all variables x and terms t . The most important property of the composition function is that applying the composition of two substitutions θ and λ to an expression e is the same as applying first one and then the other; that is, $e(\theta\lambda) = (e\theta)\lambda$. The empty substitution can be shown to be an identity under composition; that is, $\{ \}\theta = \theta\{ \} = \theta$, for all substitutions θ . Also, composition can be shown to be associative; that is, $\theta(\lambda\rho) = (\theta\lambda)\rho$ for all substitutions θ , λ , and ρ .

The definition of composition suggests a way of computing it. For example,

$$\begin{aligned} \{y \leftarrow g(z)\}\{y \leftarrow x, z \leftarrow b\} &= (y \leftarrow g(b)) \circ \{y \leftarrow x, z \leftarrow b\} \\ &= \{y \leftarrow g(b), z \leftarrow b\} \end{aligned}$$

and therefore

$$\begin{aligned} \{x \leftarrow y, y \leftarrow g(z)\}\{y \leftarrow x, z \leftarrow b\} &= (x \leftarrow x) \circ \{y \leftarrow g(b), z \leftarrow b\} \\ &= \{y \leftarrow g(b), z \leftarrow b\}. \end{aligned}$$

Note that the composition of substitutions is not commutative. For example, $\{x \leftarrow y\}\{y \leftarrow x\} = \{y \leftarrow x\}$ and $\{y \leftarrow x\}\{x \leftarrow y\} = \{x \leftarrow y\}$, but $\{y \leftarrow x\} \neq \{x \leftarrow y\}$.

A substitution θ is said to be *more general* than a substitution θ' if there exists a substitution λ such that $\theta\lambda = \theta'$. For example, the substitution $\theta : \{x \leftarrow y\}$ is more general than the substitution $\theta' : \{x \leftarrow a, y \leftarrow a\}$, because

$$\theta\{y \leftarrow a\} = \{x \leftarrow y\}\{y \leftarrow a\} = \{x \leftarrow a, y \leftarrow a\} = \theta'.$$

On the other hand, $\theta : \{x \leftarrow y\}$ is not more general than the substitution $\phi : \{x \leftarrow a\}$, because there is no substitution λ such that

$$\theta\lambda = \{x \leftarrow y\}\lambda = \{x \leftarrow a\} = \phi.$$

A substitution is regarded as more general than itself, because $\theta\{ \} = \theta$ for any substitution θ . It is possible for two distinct substitutions to be more general than each other. For example, $\theta : \{x \leftarrow y\}$ and $\theta' : \{y \leftarrow x\}$ are more general than each other, because

$$\theta\{y \leftarrow x\} = \{x \leftarrow y\}\{y \leftarrow x\} = \{y \leftarrow x\} = \theta'$$

and

$$\theta'\{x \leftarrow y\} = \{y \leftarrow x\}\{x \leftarrow y\} = \{x \leftarrow y\} = \theta.$$

UNIFIERS

A substitution θ is said to be a *unifier* of two expressions e and \tilde{e} if

$$e\theta = \tilde{e}\theta,$$

that is, if $e\theta$ and $\tilde{e}\theta$ are identical expressions. Two expressions are *unifiable* if they have a unifier.

For example, the substitution

$$\theta : \{x \leftarrow b, y \leftarrow z\}$$

is a unifier of the two expressions

$$e : f(x, z) \quad \text{and} \quad \tilde{e} : f(b, y),$$

because $e\theta = \tilde{e}\theta = f(b, z)$. Thus, e and \tilde{e} are unifiable. The substitutions

$$\phi : \{x \leftarrow b, z \leftarrow y\}$$

and

$$\rho : \{x \leftarrow b, y \leftarrow w, z \leftarrow w\}$$

are also unifiers of these two expressions. Thus, unifiers of expressions are not unique.

The expressions $p(a)$ and $p(b)$ are clearly not unifiable and neither are the expressions $q(x, f(x))$ and $q(g(y), y)$. The expressions x and $f(x)$ are also not unifiable. Because x is a proper subexpression of $f(x)$, we know $x\theta$ is a proper subexpression of $(f(x))\theta$, for any substitution θ ; hence $x\theta$ and $(f(x))\theta$ are not identical.

A substitution θ is said to be a *most-general unifier* of two expressions e and \tilde{e} if θ is a unifier of e and \tilde{e} and if θ is more general than any unifier of e and \tilde{e} . For example, the distinct substitutions $\theta : \{x \leftarrow b, y \leftarrow z\}$ and $\phi : \{x \leftarrow b, z \leftarrow y\}$ are both most general unifiers of the expressions $e : f(x, z)$ and $\tilde{e} : f(b, y)$. Thus, most-general unifiers are not unique. It is clear, however, that all most-general unifiers of two expressions are *equally general*, i.e., each is more general than any of the others.

There is a *unification algorithm* (Robinson [65]) for determining whether a given pair of expressions is unifiable and, if so, for producing a most general unifier.

We can extend the notion of unifier to apply to a list of pairs of expressions. A substitution θ is said to be a *simultaneous unifier* of the list

$$\langle\langle e_1, \tilde{e}_1 \rangle, \langle e_2, \tilde{e}_2 \rangle, \dots, \langle e_n, \tilde{e}_n \rangle\rangle$$

of pairs of expressions if

$$e_1\theta = \tilde{e}_1\theta, \quad e_2\theta = \tilde{e}_2\theta, \quad \dots, \quad \text{and} \quad e_n\theta = \tilde{e}_n\theta.$$

(Note that we do not require that $e_i\theta = e_j\theta$, for distinct i and j .) We may also say that θ is a *simultaneous unifier* of e_1 and \tilde{e}_1 , of e_2 and \tilde{e}_2 , ..., and of e_n and \tilde{e}_n . A list of pairs of expressions is *simultaneously unifiable* if it has a simultaneous unifier.

A list may fail to be simultaneously unifiable even though the expressions of each pair it contains are unifiable independently. For example, the list of pairs

$$\langle\langle x, g(y) \rangle, \langle f(x), y \rangle\rangle$$

is not simultaneously unifiable, even though the expressions x and $g(y)$ are unifiable, by the substitution $\{x \leftarrow g(y)\}$, and the expressions $f(x)$ and y are unifiable, by the substitution $\{y \leftarrow f(x)\}$.

For any list of pairs of expressions, a simultaneous unifier is *most general* if it is more general than any other simultaneous unifier.

We can extend the notion of unifier further to apply to a list of lists of expressions. A substitution θ is said to be a *simultaneous unifier* of the list

$$\langle \langle e_1, \tilde{e}_1, \tilde{\tilde{e}}_1, \dots \rangle, \langle e_2, \tilde{e}_2, \tilde{\tilde{e}}_2, \dots \rangle, \dots, \langle e_n, \tilde{e}_n, \tilde{\tilde{e}}_n, \dots \rangle \rangle,$$

of lists of expressions if

$$\begin{aligned} e_1\theta &= \tilde{e}_1\theta = \tilde{\tilde{e}}_1\theta = \dots \\ e_2\theta &= \tilde{e}_2\theta = \tilde{\tilde{e}}_2\theta = \dots \\ &\vdots \\ e_n\theta &= \tilde{e}_n\theta = \tilde{\tilde{e}}_n\theta = \dots \end{aligned}$$

We may also say that θ is a simultaneous unifier of $e_1, \tilde{e}_1, \tilde{\tilde{e}}_1, \dots$, of $e_2, \tilde{e}_2, \tilde{\tilde{e}}_2, \dots$, and of $e_n, \tilde{e}_n, \tilde{\tilde{e}}_n, \dots$. The notion of most-general simultaneous unifier and the unification algorithm may be extended accordingly. The notation is more complex but the concepts are the same.

SUBSTITUTION AND REPLACEMENT

We sometimes find it convenient to use the replacement and substitution notations together. Suppose s , t , and e are expressions, where s and t are either both sentences or both terms. Let θ be a substitution. If we write e as $e[s]$, then

$$e\theta[t]$$

denotes the expression obtained by replacing every occurrence of $s\theta$ in $e\theta$ with t . If we write e as $e\langle s \rangle$, then

$$e\theta\langle t \rangle$$

denotes the expression obtained by replacing certain (perhaps none) of the occurrences of $s\theta$ in $e\theta$ with t .

For example, consider the expression

$$e: p(f(x, a)) \text{ or } q(f(x, y)) \text{ or } r(f(b, a))$$

and the substitution

$$\theta: \{x \leftarrow b, y \leftarrow a\}.$$

If we write e as $e[f(x, a)]$, then $e\theta[g(c)]$ is

$$p(g(c)) \text{ or } q(g(c)) \text{ or } r(g(c)).$$

Note that two of the replaced occurrences of $f(x, a)\theta$ in $e\theta$ do not correspond to occurrences of $f(x, a)$ in e ; they were created by application of the substitution θ .

INTERPRETATIONS

We shall use the Herbrand notion of interpretation, in which the elements of the domain are identified with the terms of the language.

An *interpretation* I is an assignment of truth values, either T (true) or F (false), to every ground proposition (i.e., to every proposition that contains no variables). If I assigns T [or F] to a ground proposition, that proposition is said to be *true* [or *false*] *under* I . The truth [or falseness] of a nonpropositional ground sentence under an interpretation I may be determined from that of its propositional constituents by the familiar semantic rules for the logical connectives.

A nonground sentence P is *true under* I if every ground instance of P is true under I ; otherwise, P is *false under* I . Note that, according to this definition, free variables have a tacit universal quantification.

We can now define the notions of implication and equivalence between sentences. The sentences P_1, P_2, P_3, \dots *imply* a sentence Q if, for any interpretation I ,

if P_1, P_2, P_3, \dots are all true under I ,
then Q is true under I .

Note that if P implies Q , it is not necessarily the case that the sentence (*if* P *then* Q) is valid. For example, $p(x)$ implies $p(a)$, because free variables are taken to be universally quantified. But the sentence (*if* $p(x)$ *then* $p(a)$) is not valid: its instance (*if* $p(b)$ *then* $p(a)$) is false under any interpretation for which $p(b)$ is true and $p(a)$ is false.

Two sentences P and Q are *equivalent* if, for any interpretation I ,

P is true under I
if and only if
 Q is true under I .

Hence P is equivalent to Q if P implies Q and Q implies P . For example, the sentences $p(x)$ and $p(y)$ are equivalent.

Lemma (instantiation)

For any sentence \mathcal{F} and substitution θ , \mathcal{F} implies $\mathcal{F}\theta$. ┘

Both total and partial replacement exhibit the following *value* property:

Suppose P , Q , and \mathcal{F} are ground sentences and I is an interpretation. Then

if P and Q have the same truth value under I ,
then $\mathcal{F}[P]$ and $\mathcal{F}[Q]$ have the same truth value under I .

Also,

if P and Q have the same truth value under I ,
then $\mathcal{F}\langle P \rangle$ and $\mathcal{F}\langle Q \rangle$ have the same truth value under I .

A corresponding *value* property applies to multiple replacements.

Remark

The value property applies only to ground sentences, not to sentences with variables. For instance, let P be the sentence $p(x)$, let Q be the sentence *false*, and let $\mathcal{F}[P]$ be the sentence (*not* $p(x)$). Consider an interpretation I under which

$p(a)$ is true and $p(b)$ is false.

Then (by the definition of truth for a nonground sentence) $p(x)$ is false under I and hence

$p(x)$ and *false* have the same truth value under I .

On the other hand (by the definition again) $\text{not } p(x)$ is also false under I and hence

$(\text{not } p(x))$ and (not false) do not have the same truth value under I ,

contradicting the conclusion of the value property. \blacksquare

THEORIES

A *theory* is a set of sentences \mathcal{T} that is closed under logical implication: If \mathcal{T} implies a sentence \mathcal{P} then \mathcal{P} belongs to \mathcal{T} . A member of a theory \mathcal{T} is also said to be *valid* in \mathcal{T} .

A theory \mathcal{T} is said to be *defined* by a set of sentences \mathcal{A} if \mathcal{T} is precisely the set of sentences implied by \mathcal{A} . We shall also say that \mathcal{A} is a set of *axioms* for \mathcal{T} .

An interpretation I is said to be a *model* for a theory \mathcal{T} if every sentence of \mathcal{T} is true under I .

For example, let \mathcal{T} be the set of sentences implied by the transitivity axiom,

if $x \prec y$ *and* $y \prec z$
then $x \prec z$,

and the *irreflexivity* axiom,

not $x \prec x$.

Then \mathcal{T} is a theory, defined by these axioms. The *asymmetry* property

if $x \prec y$
then *not* $y \prec x$

is a (valid) sentence of this theory.

RELATIONS

We need some special terminology for speaking about relations. Henceforth, let us consider a particular theory. When we speak of validity, we shall mean validity in that theory.

Let p and q be n -ary relations. Then we say that p *implies* q if

if $p(x_1, x_2, \dots, x_n)$ *then* $q(x_1, x_2, \dots, x_n)$

is valid (in the theory under discussion). We also say that p is *equivalent* to q if

$p(x_1, x_2, \dots, x_n) \equiv q(x_1, x_2, \dots, x_n)$

is valid.

Let \prec be an arbitrary binary relation. We shall say that, over a given theory, \prec is *reflexive* if

$x \prec x$

is valid (in the theory); \prec is *irreflexive* if

not $(x \prec x)$

is valid; \prec is *total* if

$x \prec y$ *or* $x = y$ *or* $y \prec x$

is valid; \prec is *transitive* if

$$\text{if } (x \prec y \text{ and } y \prec z) \text{ then } x \prec z$$

is valid; and \prec is *symmetric* if

$$\text{if } x \prec y \text{ then } y \prec x$$

is valid.

We regard logical connectives as relations on the set of truth values $\{T, F\}$. For instance, the implication connective (*if* P *then* Q) is the relation that holds if P has value F or if P and Q both have value T ; we may read it as " P is *false* than (or as false as) Q ." The equivalence connective $P \equiv Q$ is simply the equality relation on $\{T, F\}$. Note that, viewed as binary relations, the implication connective *if-then* is reflexive, total, and transitive, and the equivalence connective \equiv is reflexive, transitive, and symmetric.

ASSOCIATED RELATIONS

For each binary relation, we shall be concerned with certain associated relations.

Consider an arbitrary binary relation $x \prec y$ (read as " x is related to y "). The *reflexive closure* \preceq of \prec is defined by

$$x \preceq y \equiv (x \prec y \text{ or } x = y).$$

The *irreflexive restriction* \prec of \prec is defined by

$$x \prec y \equiv (x \preceq y \text{ and not } (x = y)).$$

The *inverse* \succ of \prec is defined by

$$x \succ y \equiv y \prec x.$$

The *negation* \nprec of \prec is defined by

$$x \nprec y \equiv \text{not } (x \preceq y).$$

We use \succ and \preceq to denote the inverses of \prec and \preceq , respectively, and \nprec and \npreceq to denote their negations. If we are using the prefix notation $p(x, y)$ for a binary relation, we denote its reflexive closure by $\bar{p}(x, y)$, its irreflexive restriction by $\hat{p}(x, y)$, and its negation by $\neg p(x, y)$.

The following proposition connects the relations associated with a given binary relation:

Proposition (negation of associated relations)

Consider an arbitrary binary relation \prec .

The negation \nprec of the reflexive closure of \prec is equivalent to the irreflexive restriction of its negation \nprec , that is,

$$x \npreceq y \quad \text{if and only if} \quad (x \nprec y \text{ and not } (x = y)).$$

The negation \nprec of the irreflexive restriction of \prec is equivalent to the reflexive closure of its negation \nprec , that is,

$$x \nprec y \quad \text{if and only if} \quad (x \npreceq y \text{ or } x = y). \quad \blacksquare$$

3. RELATIONAL POLARITY

We are now ready to define our key notion, the polarity of a subexpression with respect to a given binary relation. We actually define the polarity of a subexpression with respect to two binary relations, \prec_1 and \prec_2 . This notion is to be defined so that, if the subexpression is positive, replacing that subexpression with a larger expression (with respect to \prec_1) will make the entire expression larger (with respect to \prec_2). Similarly, if the subexpression is negative, replacing that subexpression with a smaller expression (with respect to \prec_1) will make the entire expression larger (with respect to \prec_2).

We begin by defining polarity for the arguments of an operator (i.e., function or relation).

Definition (polarity of an operator)

Let f be an n -ary operator and \prec_1 and \prec_2 be binary relations. Then

- f is *positive* over its i th argument with respect to \prec_1 and \prec_2 if the sentence

$$\begin{array}{l} \text{if } x \prec_1 y \\ \text{then } f(z_1, \dots, z_{i-1}, x, z_{i+1}, \dots, z_n) \prec_2 f(z_1, \dots, z_{i-1}, y, z_{i+1}, \dots, z_n) \end{array}$$

is valid. In other words, replacing x with a larger element y makes

$$f(z_1, \dots, z_{i-1}, x, z_{i+1}, \dots, z_n)$$

larger.

- f is *negative* over its i th argument with respect to \prec_1 and \prec_2 if the sentence

$$\begin{array}{l} \text{if } x \prec_1 y \\ \text{then } f(z_1, \dots, z_{i-1}, y, z_{i+1}, \dots, z_n) \prec_2 f(z_1, \dots, z_{i-1}, x, z_{i+1}, \dots, z_n) \end{array}$$

is valid. In other words, replacing y with a smaller element x makes

$$f(z_1, \dots, z_{i-1}, y, z_{i+1}, \dots, z_n)$$

larger. \blacksquare

We illustrate this notion with two examples.

Example

Suppose our theory includes the finite sets and the nonnegative integers. Take $f(z)$ to be the cardinality function $\text{card}(z)$, which maps each set into the number of elements it contains. Take \prec_1 to be the subset relation \subseteq over the finite sets and \prec_2 to be the weak less-than relation \leq over the nonnegative integers.

Then the card function is positive over its first (and only) argument with respect to the relations \subseteq and \leq , because the sentence

$$\begin{array}{l} \text{if } x \subseteq y \\ \text{then } \text{card}(x) \leq \text{card}(y) \end{array}$$

is valid (in the theory). \blacksquare

Example

Consider the theory of the integers. Take $f(z_1, z_2)$ to be the less-than relation $z_1 < z_2$. Take $x \prec_1 y$ to be the predecessor relation $x \prec_{\text{pred}} y$, which holds if $x = y - 1$, and take \prec_2 to be the *if-then* connective. (Recall that we regard connectives as relations on the set of truth values.)

Then the less-than relation $<$ is negative over its first argument with respect to \prec_{pred} and *if-then*, because the sentence

$$\begin{array}{l} \text{if } x \prec_{pred} y \\ \text{then if } y < z_2 \text{ then } x < z_2 \end{array}$$

is valid. Also, $<$ is positive over its second argument with respect to \prec_{pred} and *if-then*, because the sentence

$$\begin{array}{l} \text{if } x \prec_{pred} y \\ \text{then if } z_1 < x \text{ then } z_1 < y \end{array}$$

is valid. \perp

It follows from the definition that, for any n -ary operator f and binary relations \prec_1 and \prec_2 ,

- f is positive over its i th argument with respect to \prec_1 and \prec_2
- if and only if
- f is negative over its i th argument with respect to \succ_1 and \prec_2
- if and only if
- f is negative over its i th argument with respect to \prec_1 and \succ_2
- if and only if
- f is positive over its i th argument with respect to \succ_1 and \succ_2 .

When we say that a relation $p(z_1, \dots, z_n)$ is positive or negative over its i th argument with respect to a single relation \prec_1 , without mentioning a second relation \prec_2 , we shall by convention take \prec_2 to be the *if-then* connective. Thus in the above example we may simply say that $<$ is negative over its first argument and positive over its second argument, with respect to \prec_{pred} .

Every relation is both positive and negative over each of its arguments with respect to the equality relation $=$, because the sentences

$$\begin{array}{l} \text{if } x = y \\ \text{then if } p(z_1, \dots, x, \dots, z_n) \\ \text{then } p(z_1, \dots, y, \dots, z_n) \end{array} \quad \text{and} \quad \begin{array}{l} \text{if } x = y \\ \text{then if } p(z_1, \dots, y, \dots, z_n) \\ \text{then } p(z_1, \dots, x, \dots, z_n) \end{array}$$

are valid. This is equivalent to the relational-substitutivity property of equality. Also, every function is both positive and negative over each of its arguments with respect to $=$ and $=$, because the sentences

$$\begin{array}{l} \text{if } x = y \\ \text{then } f(z_1, \dots, x, \dots, z_n) = f(z_1, \dots, y, \dots, z_n) \end{array} \quad \text{and} \quad \begin{array}{l} \text{if } x = y \\ \text{then } f(z_1, \dots, y, \dots, z_n) = f(z_1, \dots, x, \dots, z_n) \end{array}$$

are valid. This is equivalent to the functional-substitutivity property of equality.

Every connective is both positive and negative over all its arguments with respect to \equiv . For example, the *not* connective is both positive and negative over its argument with respect to \equiv , because both sentences

$$\begin{array}{l} \text{if } x \equiv y \\ \text{then if } (\text{not } x) \text{ then } (\text{not } y) \end{array} \quad \text{and} \quad \begin{array}{l} \text{if } x \equiv y \\ \text{then if } (\text{not } y) \text{ then } (\text{not } x) \end{array}$$

are valid.

When we say that a connective is positive or negative over its i th argument, without mentioning any relations \prec_1 and \prec_2 at all, we shall by convention take both \prec_1 and \prec_2 to be the *if-then* connective. Polarity in this sense is close to its ordinary use in logic. The negation connective *not* is negative in its first (and only) argument, because the sentence

$$\begin{array}{l} \text{if if } x \text{ then } y \\ \text{then if } (\text{not } y) \text{ then } (\text{not } x) \end{array}$$

is valid. The conjunction connective *and* and the disjunction connective *or* are positive over both their arguments. The implication connective *if-then* is negative in its first argument, but positive in its second. The equivalence connective \equiv has no polarity in either argument. The conditional connective *if-then-else* has no polarity in its first argument, but is positive in its second and third argument.

Note that a binary relation \prec is transitive if and only if it is negative with respect to \prec itself over its first argument, because the polarity condition

$$\begin{array}{l} \text{if } x \prec y \\ \text{then if } y \prec z \text{ then } x \prec z \end{array}$$

is equivalent to the definition of transitivity. Also, \prec is transitive if and only if it is positive with respect to \prec over its second argument.

We are now ready to define polarity for the subexpressions of a given expression. The definition is inductive.

Definition (polarity of a subexpression)

Let \prec_1 and \prec_2 be binary relations. Then

- An expression s is *positive in s itself* with respect to \prec_1 and \prec_2 if \prec_1 implies \prec_2 .
- An expression s is *negative in s itself* with respect to \prec_1 and \prec_2 if \prec_1 implies \succ_2 .

Let f be an n -ary operator and e_1, e_2, \dots, e_n be expressions. Consider an occurrence of s in one of the expressions e_i . Then

- The occurrence of s is *positive in $f(e_1, e_2, \dots, e_n)$* with respect to \prec_1 and \prec_2 if there exists a binary relation \prec such that
 - the polarity of the occurrence of s in e_i with respect to \prec_1 and \prec is the same as
 - the polarity of f over its i th argument with respect to \prec and \prec_2 .
- The occurrence of s is *negative in $f(e_1, e_2, \dots, e_n)$* with respect to \prec_1 and \prec_2 if there exists a binary relation \prec such that
 - the polarity of the occurrence of s in e_i with respect to \prec_1 and \prec is opposite to
 - the polarity of f over its i th argument with respect to \prec and \prec_2 .

Furthermore, if f has no polarity over its i th argument or if s has no polarity in e_i , then s has no polarity in $f(e_1, e_2, \dots, e_n)$. On the other hand, if s has both polarities in e_i and f has some polarity over its i th argument, or if f has both polarities over its i th argument and s has some polarity in e_i , then s automatically has both polarities in $f(e_1, e_2, \dots, e_n)$. \blacksquare

Remark

For any binary relation \prec , any expression s is positive in s itself with respect to \prec and \prec (because \prec implies \prec). Similarly, s is negative in s with respect to \prec and \succ .

If f is positive over its i th argument with respect to \prec_1 and \prec_2 , then, for any expressions e_1, e_2, \dots, e_n , the occurrence of e_i is positive in $f(e_1, \dots, e_i, \dots, e_n)$ with respect to \prec_1 and \prec_2 . For take \prec to be \prec_1 . Then the polarity of e_i in e_i itself is positive with respect to \prec_1 and \prec_1 . Also, f is positive over its i th argument with respect to \prec_1 and \prec_2 . Because these two polarities are the same, e_i is positive in $f(e_1, \dots, e_i, \dots, e_n)$ with respect to \prec_1 and \prec_2 .

Similarly, if f is negative over its i th argument, then e_i is negative in $f(e_1, \dots, e_i, \dots, e_n)$, with respect to \prec_1 and \prec_2 . ┘

We may indicate the polarity of a subexpression s by annotating it s^+ , s^- , or s^\pm .

For example, suppose our theory includes the theories of sets and nonnegative integers. The occurrence of s in the sentence

$$\text{card}(s^-) < m$$

is negative with respect to the subset relation \subseteq and the *if-then* connective. For note that card is positive over its argument with respect to \subseteq and \leq and that $<$ is negative over its first argument with respect to \leq and *if-then*. Therefore, by our remark, we know that s is positive in $\text{card}(s)$ with respect to \subseteq and \leq and that $\text{card}(s)$ is negative in $\text{card}(s) < m$ with respect to \leq and *if-then*. By the definition, taking \prec_1 to be \subseteq , \prec to be \leq , and \prec_2 to be *if-then*, we conclude that s is negative in $\text{card}(s) < m$ with respect to \subseteq and *if-then*.

When we say that an occurrence of a subexpression is positive or negative in a sentence with respect to a single relation \prec_1 , without mentioning a second relation \prec_2 , we shall again take \prec_2 to be the *if-then* connective. When we say that an occurrence of a subsentence is positive or negative in a sentence, without mentioning any relation at all, we shall again take both \prec_1 and \prec_2 to be *if-then*.

It can be established from the definition that, for expressions s and t and binary relations \prec_1 and \prec_2 ,

- an occurrence of s is positive in t with respect to \prec_1 and \prec_2
- if and only if
- the occurrence of s is negative in t with respect to \succ_1 and \prec_2
- if and only if
- the occurrence of s is negative in t with respect to \prec_1 and \succ_2
- if and only if
- the occurrence of s is positive in t with respect to \succ_1 and \succ_2 .

This is analogous to our previous result concerning polarity for the argument of an operator.

Suppose an occurrence of s is positive [or negative] in t with respect to \prec_1 and \prec_2 . Then if $\tilde{\prec}_1$ is a binary relation that implies \prec_1 , then s is positive [or negative, respectively] in t with respect to $\tilde{\prec}_1$ and \prec_2 . Similarly, if \prec_2 implies a binary relation $\tilde{\prec}_2$, then s is positive [or negative, respectively] in t with respect to \prec_1 and $\tilde{\prec}_2$.

We can establish the following result:

Lemma (polarity operator)

Let \prec_1 and \prec_2 be binary relations, f be an n -ary operator, and e_1, e_2, \dots, e_n be expressions. Consider an occurrence of s in one of the expressions e_i such that s has some polarity in $f(e_1, e_2, \dots, e_n)$ with respect to \prec_1 and \prec_2 .

Then there exists a binary relation \prec such that

$$f \text{ is positive over its } i\text{th argument with respect to } \prec \text{ and } \prec_2$$

and

the polarity of the occurrence of s in $f(e_1, e_2, \dots, e_n)$ with respect to \prec_1 and \prec_2
is the same as
the polarity of the occurrence of s in e_i with respect to \prec_1 and \prec . ┘

Proof

Consider the case in which the occurrence of s is positive in $f(e_1, e_2, \dots, e_n)$ with respect to \prec_1 and \prec_2 . According to the definition, this means that there exists a binary relation $\tilde{\prec}$ such that

the polarity of the occurrence of s in e_i with respect to \prec_1 and $\tilde{\prec}$
is the same as
the polarity of f over its i th argument with respect to $\tilde{\prec}$ and \prec_2 .

If f is positive over its i th argument with respect to $\tilde{\prec}$ and \prec_2 , then the occurrence of s is positive in e_i with respect to \prec_1 and $\tilde{\prec}$, and we can simply take \prec to be $\tilde{\prec}$.

On the other hand, if f is negative over its i th argument with respect to $\tilde{\prec}$ and \prec_2 , then the occurrence of s is negative in e_i with respect to \prec_1 and $\tilde{\prec}$. By previous remarks, this means that f is positive over its i th argument with respect to the inverse relation $\tilde{\prec}^{-1}$ and \prec_2 , and the occurrence of s is positive in e_i with respect to \prec_1 and the inverse relation $\tilde{\prec}^{-1}$. Hence we can take \prec to be $\tilde{\prec}^{-1}$.

The case in which s is negative in $f(e_1, e_2, \dots, e_n)$ is treated similarly. ┘

Polarities of subexpressions of subexpressions can be composed according to the following result.

Lemma (polarity composition)

Consider an occurrence of a subexpression r in an expression s and an occurrence of s in an expression t . Then the polarity of the occurrence of r is positive [or negative] in t with respect to binary relations \prec_1 and \prec_2 if and only if there exists a binary relation \prec such that

the polarity of the occurrence of r in s with respect to \prec_1 and \prec
is the same as [or opposite to, respectively]
the polarity of the occurrence of s in t with respect to \prec and \prec_2 . ┘

For instance, if r is negative in s and s is negative in t then r is positive in t , with respect to the appropriate binary relations. If r has both polarities in s and s has some polarity in t , then r has both polarities in t .

We can now establish the fundamental property of polarity.

Lemma (polarity replacement)

For any binary relations \prec_1 and \prec_2 and expression $e\langle x^+, y^- \rangle$, the sentence

if $x \prec_1 y$
then $e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^1$

is valid. Here $e\langle y^+, x^- \rangle^1$ is the result of replacing in $e\langle x^+, y^- \rangle$ precisely one positive occurrence of x with y or negative occurrence of y with x (we assume that such an occurrence exists) where the polarity is taken in $e\langle x^+, y^- \rangle$ with respect to \prec_1 and \prec_2 . ┘

Example

Suppose our theory includes the theories of lists and nonnegative integers. Take \prec_1 to be the tail relation $x \prec_{tail} y$, which is true if

not ($y = []$) and $x = tail(y)$,

that is, if y is nonempty and x is the list of all but the first element of y . Take \prec_2 to be the predecessor relation \prec_{pred} . Take $e\langle x^+, y^- \rangle$ to be the expression

$$length(x^+) + length(x^-),$$

where the function $length(x)$ yields the number of elements in the list x .

Note that each occurrence of x is positive in $length(x) + length(x)$ with respect to \prec_{tail} and \prec_{pred} , as indicated by the annotations. For, each occurrence is positive in $length(x)$ with respect to \prec_{tail} and \prec_{pred} , and the plus function $+$ is positive over either of its arguments with respect to \prec_{pred} and \prec_{pred} .

Therefore, according to the lemma, the sentence

$$\begin{array}{l} \text{if } x \prec_{tail} y \\ \text{then } length(x) + length(x) \prec_{pred} length(y) + length(x) \end{array}$$

is valid, because $length(y) + length(x)$ is the result of replacing one positive occurrence of x in $length(x) + length(x)$ with y . Also, according to the lemma, the sentence

$$\begin{array}{l} \text{if } x \prec_{tail} y \\ \text{then } length(x) + length(x) \prec_{pred} length(x) + length(y) \end{array}$$

is valid, because $length(x) + length(y)$ is the result of replacing one positive occurrence of x in $length(x) + length(x)$ with y .

On the other hand, the lemma does not allow us to conclude that

$$\begin{array}{l} \text{if } x \prec_{tail} y \\ \text{then } length(x) + length(x) \prec_{pred} length(y) + length(y) \end{array}$$

is valid, because $length(y) + length(y)$ is obtained by replacing two, not one, positive occurrences of x in $length(x) + length(x)$ with y . In fact, this third sentence is not valid. \blacksquare

We now prove the lemma.

Proof (polarity replacement lemma)

For any arbitrary binary relation \prec_1 , suppose that

$$x \prec_1 y.$$

We show that, for any expression $e\langle x^+, y^- \rangle$, we have, for any binary relation \prec_2 ,

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^1.$$

The proof is by induction on the structure of $e\langle x^+, y^- \rangle$. In other words, we show the desired conclusion for an arbitrary expression $e\langle x^+, y^- \rangle$, under the induction hypothesis that, for any proper subexpression $\tilde{e}\langle x^+, y^- \rangle$ of $e\langle x^+, y^- \rangle$, we have, for any binary relation $\tilde{\prec}_2$,

$$\tilde{e}\langle x^+, y^- \rangle \tilde{\prec}_2 \tilde{e}\langle y^+, x^- \rangle^1.$$

As in the statement of the lemma, $\tilde{e}\langle y^+, x^- \rangle^1$ is obtained from $\tilde{e}\langle x^+, y^- \rangle$ by replacing precisely one occurrence of x or y , of suitable polarity with respect to \prec_1 and $\tilde{\prec}_2$.

The proof distinguishes among several subcases.

Case: The expression $e\langle x^+, y^- \rangle$ is simply x

Then, because the replaced variable x is positive in x , with respect to \prec_1 and \prec_2 , we have (by the definition of polarity) that \prec_1 implies \prec_2 .

In this case, $e\langle y^+, x^- \rangle^1$ is y , and we must show

$$x \prec_2 y.$$

But this follows from our supposition that $x \prec_1 y$, because \prec_1 implies \prec_2 .

Case: The expression $e\langle x^+, y^- \rangle$ is simply y

Then, because the replaced variable y is negative in y with respect to \prec_1 and \prec_2 , we have (by the definition of polarity) that \prec_1 implies \succ_2 .

In this case, $e\langle y^+, x^- \rangle$ is x , and we must show that

$$y \prec_2 x,$$

or, equivalently, that

$$x \succ_2 y.$$

But this follows from our supposition that $x \prec_1 y$, because \prec_1 implies \succ_2 .

Case: $e\langle x^+, y^- \rangle$ is of form $f(e_1, e_2, \dots, e_n)$, where f is an n -ary operator

The replaced occurrence of x [or y] must occur in one of the arguments e_i of f . Because this occurrence is positive [or negative, respectively] in $f(e_1, e_2, \dots, e_n)$ with respect to \prec_1 and \prec_2 , we know (by the *polarity operator lemma*) that there exists a binary relation \prec such that

f is positive over its i th argument with respect to \prec and \prec_2

and

the polarity of the replaced occurrence of x [or y] in e_i with respect to \prec_1 and \prec is the same as the polarity of the replaced occurrence of x [or y] in $f(e_1, e_2, \dots, e_n)$, that is, $e\langle x^+, y^- \rangle$, with respect to \prec_1 and \prec_2 .

Let us therefore write e_i as $e_i\langle x^+, y^- \rangle$.

Because $e_i\langle x^+, y^- \rangle$ is a proper subexpression of $e\langle x^+, y^- \rangle$, we can apply our induction hypothesis, taking $\tilde{e}\langle x^+, y^- \rangle$ to be $e_i\langle x^+, y^- \rangle$ and $\tilde{\prec}_2$ to be \prec , to conclude that

$$e_i\langle x^+, y^- \rangle \prec e_i\langle y^+, x^- \rangle^1.$$

Therefore (by the definition of polarity of an operator, because f is positive over its i th argument with respect to \prec and \prec_2), we have

$$f(e_1, \dots, e_i\langle x^+, y^- \rangle, \dots, e_n) \prec_2 f(e_1, \dots, e_i\langle y^+, x^- \rangle, \dots, e_n),$$

that is,

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^1,$$

as we wanted to show. This completes the proof. \blacksquare

The polarity replacement lemma allows us to replace precisely one occurrence of a variable. If we know more about the relation \prec_2 , we can establish stronger versions of the lemma. In particular, if we know that \prec_2 is transitive, we can replace one or more occurrences of the variable.

Lemma (transitive polarity replacement)

For any binary relations \prec_1 and \prec_2 and expression $e\langle x^+, y^- \rangle$, where \prec_2 is transitive, the sentence

$$\begin{array}{l} \text{if } x \prec_1 y \\ \text{then } e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^n \end{array}$$

is valid for every positive integer n . Here $e\langle y^+, x^- \rangle^n$ is the result of replacing in $e\langle x^+, y^- \rangle$ precisely n positive occurrences of x with y or negative occurrences of y with x , where the polarity is taken in $e\langle x^+, y^- \rangle$ with respect to \prec_1 and \prec_2 . \blacksquare

Note that we can replace occurrences of both x and y in the same expression; precisely n replacements are made altogether. Also, the lemma requires that at least one replacement be made.

Example

Suppose our theory includes the theories of both lists and integers. Take $e\langle x^+, y^- \rangle$ to be the expression

$$e\langle x^+, y^- \rangle : \text{length}(x^+) + (\text{length}(x^+) - \text{length}(y^-)).$$

Take \prec_1 to be the tail relation \prec_{tail} (defined in a previous example) and \prec_2 to be the less-than relation $<$. Note that, with respect to \prec_{tail} and $<$, both occurrences of x are positive and the occurrence of y is negative in $e\langle x^+, y^- \rangle$; also $<$ is transitive. According to the lemma, the following sentences (among others) are valid: the sentence

$$\begin{array}{l} \text{if } x \prec_{tail} y \\ \text{then } \text{length}(x) + (\text{length}(x) - \text{length}(y)) < \text{length}(y) + (\text{length}(y) - \text{length}(y)), \end{array}$$

for which both occurrences of x in $e\langle x^+, y^- \rangle$ have been replaced, and

$$\begin{array}{l} \text{if } x \prec_{tail} y \\ \text{then } \text{length}(x) + (\text{length}(x) - \text{length}(y)) < \text{length}(x) + (\text{length}(y) - \text{length}(x)), \end{array}$$

for which one occurrence of x and one of y in $e\langle x^+, y^- \rangle$ have been replaced.

On the other hand, the lemma does not allow us to conclude that

$$\begin{array}{l} \text{if } x \prec_{tail} y \\ \text{then } \text{length}(x) + (\text{length}(x) - \text{length}(y)) < \text{length}(x) + (\text{length}(x) - \text{length}(y)), \end{array}$$

is valid, because no replacements of x or of y in $e\langle x^+, y^- \rangle$ have been replaced. In fact, this final sentence is clearly not valid. \blacksquare

We now prove the lemma

Proof (transitive polarity replacement lemma)

We assume throughout that polarity is with respect to \prec_1 and \prec_2 . We suppose that

$$x \prec_1 y$$

and show that

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^n,$$

for every positive integer n . The proof is by induction on n .

Base Case: $n = 1$.

In this case, precisely one replacement is made. The desired result

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^1$$

follows from the original polarity replacement lemma.

Inductive Step:

For an arbitrary positive integer k , we assume inductively that

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^k$$

and show that

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^{k+1}.$$

Observe that $e\langle y^+, x^- \rangle^{k+1}$ can be obtained from $e\langle y^+, x^- \rangle^k$ by replacing precisely one positive occurrence of x with y or one negative occurrence of y with x . Therefore, by the original polarity replacement lemma, we have

$$e\langle y^+, x^- \rangle^k \prec_2 e\langle y^+, x^- \rangle^{k+1}.$$

Because our induction hypothesis is that $e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^k$, and because we have assumed that \prec_2 is transitive, we can conclude that

$$e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle^{k+1},$$

as we wanted to show. \blacksquare

If \prec_2 is transitive, the above lemma allows us to replace one or more occurrences of a variable. If \prec_2 is both reflexive and transitive, the following lemma allows us to replace zero, one, or more occurrences.

Lemma (reflexive transitive polarity replacement)

For any binary relations \prec_1 and \prec_2 and expression $e\langle x^+, y^- \rangle$, where \prec_2 is both reflexive and transitive, the sentence

$$\begin{array}{l} \text{if } x \prec_1 y \\ \text{then } e\langle x^+, y^- \rangle \prec_2 e\langle y^+, x^- \rangle \end{array}$$

is valid. Here $e\langle y^+, x^- \rangle$ is the result of replacing in $e\langle x^+, y^- \rangle$ certain positive occurrences of x with y and certain negative occurrences of y with x , where polarity is taken in $e\langle x^+, y^- \rangle$ with respect to \prec_1 and \prec_2 . \blacksquare

This lemma, as opposed to the *transitive polarity replacement* lemma, admits the possibility of replacing no occurrences at all of x or y in $e\langle x^+, y^- \rangle$.

Example

Suppose our theory includes the theories of both finite sets and integers. Take $e\langle x^+, y^- \rangle$ to be the expression

$$e\langle x^+, y^- \rangle : \text{card}(x^+ \sim y^-) - \text{card}(y^- \sim x^+)$$

where $x \sim y$ is the difference between the sets x and y , that is, the set of elements of x that do not belong to y . Take \prec_1 to be the subset relation \subseteq and \prec_2 to be the weak less-than relation \leq . Note that, with respect to \subseteq and \leq , both occurrences of x are positive and both occurrences of y are negative in $e\langle x^+, y^- \rangle$, as the annotations indicate. Also, \leq is both transitive and reflexive.

Therefore, according to the lemma, the following sentences are valid: the sentence

$$\begin{array}{l} \text{if } x \subseteq y \\ \text{then } \text{card}(x \sim y) - \text{card}(y \sim x) \leq \text{card}(y \sim x) - \text{card}(x \sim y), \end{array}$$

for which all occurrences of x and y in $e\langle x^+, y^- \rangle$ have been replaced, and the sentence

$$\begin{array}{l} \text{if } x \subseteq y \\ \text{then } \text{card}(x \sim y) - \text{card}(y \sim x) \leq \text{card}(x \sim y) - \text{card}(y \sim x), \end{array}$$

for which no occurrences of x and y in $e\langle x^+, y^- \rangle$ have been replaced. Of course, other valid sentences can be obtained by replacing some, but not all, of the occurrences of x and y in $e\langle x^+, y^- \rangle$. \blacksquare

The proof is straightforward.

Proof (reflexive transitive polarity-replacement lemma)

In the case in which no replacements are made, $e\langle y^+, x^- \rangle$ is identical to $e\langle x^+, y^- \rangle$, and the desired result holds because we have supposed that \prec_2 is reflexive. In the case in which one or more replacements are made, the desired result follows from the *transitive polarity replacement* lemma, because we have supposed that \prec_2 is also transitive. \blacksquare

The following consequence of the polarity replacement lemma will be used most frequently:

Proposition (polarity replacement)

For any binary relation \prec and sentence $\mathcal{P}\langle x^+, y^- \rangle$, the sentence

$$\begin{array}{l} \text{if } x \prec y \\ \text{then if } \mathcal{P}\langle x^+, y^- \rangle \\ \text{then } \mathcal{P}\langle y^+, x^- \rangle \end{array}$$

is valid. Here $\mathcal{P}\langle y^+, x^- \rangle$ is the result of replacing in $\mathcal{P}\langle x^+, y^- \rangle$ certain positive occurrences of x with y and certain negative occurrences of y with x , where polarity is taken in $\mathcal{P}\langle x^+, y^- \rangle$ with respect to \prec . \blacksquare

Recall that, when we refer to polarity in a sentence with respect to a single relation \prec , we mean polarity with respect to \prec and the *if-then* connective. The proposition allows us to replace occurrences of both x and y in the same sentence and (trivially) admits the possibility that no replacements are made.

The proof is immediate.

Proof

Regarded as a relation, the *if-then* connective is reflexive and transitive. The replaced occurrences of x and y are respectively positive and negative in $\mathcal{P}\langle x^+, y^- \rangle$ with respect to \prec and *if-then*. Therefore the proposition is simply an instance of the *reflexive transitive polarity replacement* lemma, taking \prec_1 to be \prec , \prec_2 to be *if-then*, and $e\langle x^+, y^- \rangle$ to be $\mathcal{P}\langle x^+, y^- \rangle$. \blacksquare

Example

Suppose our theory includes the theories of finite sets and integers. Take $\mathcal{P}\langle x^+, y^- \rangle$ to be the sentence

$$\mathcal{P}\langle x^+, y^- \rangle: \quad a < \text{card}(x^+ \sim y^-) \text{ and } \text{card}(y^- \sim x^+) \leq b.$$

Take \prec to be the subset relation \subseteq . Note that, with respect to \subseteq , both occurrences of x are positive and both occurrences of y are negative in $\mathcal{P}\langle x^+, y^- \rangle$, as indicated by the annotations. Therefore, according to the proposition, the following sentences are valid: the sentence

$$\begin{array}{l} \text{if } x \subseteq y \\ \text{then if } a < \text{card}(x \sim y) \text{ and } \text{card}(y \sim x) \leq b \\ \text{then } a < \text{card}(x \sim x) \text{ and } \text{card}(y \sim y) \leq b, \end{array}$$

for which one occurrence of x and one occurrence of y in $\mathcal{P}\langle x^+, y^- \rangle$ has been replaced, the sentence

if $x \subseteq y$
 then if $a < \text{card}(x \sim y)$ and $\text{card}(y \sim x) \leq b$
 then $a < \text{card}(y \sim y)$ and $\text{card}(y \sim y) \leq b$,

for which both occurrences of x in $\mathcal{P}\langle x^+, y^- \rangle$ have been replaced, and the sentence

if $x \subseteq y$
 then if $a < \text{card}(x \sim y)$ and $\text{card}(y \sim x) \leq b$
 then $a < \text{card}(y \sim x)$ and $\text{card}(x \sim y) \leq b$,

for which both occurrences of x and both occurrences of y in $\mathcal{P}\langle x^+, y^- \rangle$ have been replaced. ┘

We have now developed the mathematical results on relational polarity we need in order to introduce the special-relations rules. But first, we introduce briskly our basic nonclausal deduction system.

4. NONCLAUSAL DEDUCTION

In this section we present a basic nonclausal deduction system, without any special-relations rules. This system bears some resemblance to those of Murray [82] and Stickel [82]; it is based on the system of Manna and Waldinger [80], but is simplified in several respects:

- The system presented here is a refutation system; it attempts to show that a given set of sentences is unsatisfiable. (The original system operates on a tableau of assertions and goals, and attempts to show that at least one of the goals follows from the assertions.)
- The system is presented with no program synthesis capabilities.
- The mathematical induction principle is omitted.

These simplifications have been made for purely expository purposes: the special-relations rules are compatible with a tableau theorem prover and with the induction principle and are of great use in program synthesis, our primary application.

THE DEDUCED SET

The deduction system we describe operates on a set, called the *deduced set*, of sentences in quantifier-free first-order logic. We attempt to show that a given deduced set is unsatisfiable, i.e., that there is no interpretation under which all the sentences are true.

Theorem proving in a first-order axiomatic theory can be reduced to showing the unsatisfiability of such a set. In particular, to show that a sentence \mathcal{F} is valid in a theory whose axioms are $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$, we can

- Remove the quantifiers of the sentences $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$, and *not* \mathcal{F} , by skolemization (see, for example, Chang and Lee [73], Loveland [78], or Robinson [79]).
- Show the unsatisfiability of the resulting set of quantifier-free sentences.

We do not require that the sentences be in clausal form; indeed, they can use the full set of connectives of propositional logic, including equivalence (\equiv) and the conditional (*if-then-else*).

Example

Consider the theory of the *strict partial ordering* \prec , defined by the *transitivity* axiom

$$(\forall x)(\forall y)(\forall z) \left[\begin{array}{l} \text{if } x \prec y \text{ and } y \prec z \\ \text{then } x \prec z \end{array} \right]$$

and the *irreflexivity* axiom

$$(\forall x)[\text{not } (x \prec x)].$$

Suppose we would like to show that in this theory the *asymmetry* property

$$(\forall u)(\forall v) \left[\begin{array}{l} \text{if } u \prec v \\ \text{then not } v \prec u \end{array} \right]$$

is valid. It suffices to show that the set of quantifier-free sentences

$$\begin{array}{lll} \text{if } x \prec y \text{ and } y \prec z & \text{not } (x \prec x) & \text{not } \left[\begin{array}{l} \text{if } a \prec b \\ \text{then not } (b \prec a) \end{array} \right] \\ \text{then } x \prec z & & \end{array}$$

is unsatisfiable. \perp

If the truth symbol *false* belongs to the deduced set, the set is automatically unsatisfiable, because the sentence *false* is not true under any interpretation.

Because the variables of the sentences in the deduced set are tacitly quantified universally, we can systematically rename them without changing the unsatisfiability of the set; that is, the set is unsatisfiable before the renaming if and only if it is unsatisfiable afterwards. Of course, we must replace every occurrence of a variable in the sentence with the new variable, and we must be careful not to replace distinct variables in the sentence with the same variable. The variables of the sentences in the deduced set may therefore be *standardized apart*; in other words, we may rename the variables of the sentences so that no two of them have variables in common.

For any sentence \mathcal{F} in the deduced set and any substitution θ , we may add to the set the *instance* $\mathcal{F}\theta$ of \mathcal{F} , without changing the unsatisfiability of the set. In particular, if the deduced set is unsatisfiable after the addition of the new sentence, it was also unsatisfiable before. Note that in adding the new sentence $\mathcal{F}\theta$, we do not remove the original sentence \mathcal{F} .

THE DEDUCTIVE PROCESS

In the deductive system we apply *deduction rules*, which add new sentences to the deduced set without changing its unsatisfiability. Deduction rules are expressed as follows:

$$\frac{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_m}{\mathcal{F}}$$

This means that, if the *given* sentences $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_m$ belong to the deduced set, the *conclusion* \mathcal{F} may be added. Such a rule is said to be *sound* if the given sentences $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_m$ imply the sentence \mathcal{F} . If a deductive rule is sound, its application will preserve the unsatisfiability of the deduced set.

The deductive process terminates successfully if we introduce the truth symbol *false* into the deduced set. Because deduction rules preserve unsatisfiability, and because a set of sentences containing *false* is automatically unsatisfiable, this will imply that the original deduced set was also unsatisfiable.

We include two classes of deduction rules in the basic system:

- The *transformation* rules, which replace subsentences with equivalent sentences.
- The *resolution* rule, which performs a case analysis on the truth of matching subsentences.

These rules are described in this section. In later sections, we augment the basic system with two new classes of rules:

- The *replacement* rules, which replace subexpressions with other expressions (not necessarily equivalent or equal).
- The *matching* rules, which introduce new conditions to be proved that enable subexpressions to be matched.

We first describe the transformation rules.

TRANSFORMATION RULES

The transformation rules replace subsentences of the sentences of our deduced set with propositionally equivalent, simpler sentences. For instance, the transformation rule

$$P \text{ and } true \rightarrow P$$

replaces a subsentence of form $(P \text{ and } true)$ with the corresponding sentence of form P . The simplified sentence is then added to the deduced set. (Logically speaking, the original sentence remains in the deduced set too, but, for efficiency of implementation, the original sentence need not be retained.)

We include a full set of such *true-false* transformation rules; e.g.,

$$not \ true \rightarrow false$$

$$P \text{ or } true \rightarrow true$$

$$if \ P \text{ then } false \rightarrow not \ P.$$

These rules can eliminate from a sentence any occurrence of the truth symbols *true* and *false* as a proper subsentence.

We also include such propositional simplification rules as

$$P \text{ and } P \rightarrow P$$

$$not \ not \ P \rightarrow P.$$

These rules are not logically necessary, but are included for cosmetic purposes.

The soundness of the transformation rules is evident, because each produces a sentence equivalent to the one to which it is applied.

Example

Suppose our deduced set contains the sentence

$$\begin{array}{l} \mathcal{F}: \quad \text{if } q(a) \text{ then } false \\ \quad \quad \text{or} \\ \quad \quad (not \ true) \text{ or } (not \ q(a)). \end{array}$$

(We omit parentheses when the structure of the sentence can be indicated by indenting.) This can be transformed, by application of the rule

$$if \ P \text{ then } false \rightarrow not \ P,$$

into the sentence

$$\begin{array}{c} \text{not } q(a) \\ \text{or} \\ (\text{not true}) \text{ or } (\text{not } q(a)), \end{array}$$

which may then be added to the deduced set.

The new sentence can be transformed in turn, by successive application of the rules

$$\begin{array}{l} \text{not true} \rightarrow \text{false} \\ \text{false or } \mathcal{P} \rightarrow \mathcal{P}, \\ \mathcal{P} \text{ or } \mathcal{P} \rightarrow \mathcal{P}, \end{array}$$

into the sentence

$$\text{not } q(a).$$

We shall say that the original sentence \mathcal{F} *reduces to* $(\text{not } q(a))$ under transformation. \blacksquare

Our original system (Manna and Waldinger [80]) included many more transformation rules; also, their operation was more complex. In this system, the role of these more complex rules has been assumed by the replacement rule of Section 5.

RESOLUTION RULE: GROUND VERSION

The resolution rule applies to two sentences of our set, and performs a case analysis on the truth of a common subsentence. Instances of the sentences can be formed, if necessary, to create a common subsentence; however, we first present the *ground version* of the rule, which does not form instances of these sentences.

Rule (resolution, ground version)

For any ground sentences \mathcal{P} , $\mathcal{F}[\mathcal{P}]$, and $\mathcal{G}[\mathcal{P}]$, we have

$$\frac{\begin{array}{c} \mathcal{F}[\mathcal{P}] \\ \mathcal{G}[\mathcal{P}] \end{array}}{\mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}]} \blacksquare$$

In other words, if $\mathcal{F}[\mathcal{P}]$ and $\mathcal{G}[\mathcal{P}]$ are sentences in our deduced set with a common subsentence \mathcal{P} , we can add to the set the sentence $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}])$ obtained by replacing every occurrence of \mathcal{P} in $\mathcal{F}[\mathcal{P}]$ with *false*, replacing every occurrence of \mathcal{P} in $\mathcal{G}[\mathcal{P}]$ with *true*, and taking the disjunction of the results. We shall assume that $\mathcal{F}[\mathcal{P}]$ and $\mathcal{G}[\mathcal{P}]$ have at least one occurrence each of the subsentence \mathcal{P} . We do not require that $\mathcal{F}[\mathcal{P}]$ and $\mathcal{G}[\mathcal{P}]$ be distinct sentences.

Because the resolution rule introduces new occurrences of the truth symbols *true* and *false*, it is always possible to simplify the resulting sentence immediately afterwards by application of the appropriate *true-false* rules. These subsequent transformations will sometimes be regarded as part of the resolution rule itself.

Example

Suppose our deduced set contains the sentences

$$\mathcal{F}: \text{ if } q(a) \text{ then } \boxed{p(a, b)}$$

and

$$\mathcal{G}: (\text{not } \boxed{p(a, b)}) \text{ or } (\text{not } q(a)).$$

These sentences have a common subsentence $p(a, b)$, indicated by the surrounding boxes. By application of the resolution rule, we may replace every occurrence of $p(a, b)$ in \mathcal{F} with *false*, replace every occurrence of $p(a, b)$ in \mathcal{G} with *true*, and take the disjunction of the result, obtaining the sentence

$$\begin{array}{l} \text{if } q(a) \text{ then false} \\ \text{or} \\ (\text{not true}) \text{ or } (\text{not } q(a)), \end{array}$$

which (as we have seen in a previous example) reduces under transformation to

$$\text{not } q(a).$$

This sentence may be added to the deduced set. \blacksquare

Let us show that the resolution rule is sound, and hence that it preserves the unsatisfiability of the deduced set.

Justification (resolution rule, ground version)

We must show that the given sentences $\mathcal{F}[\mathcal{P}]$ and $\mathcal{G}[\mathcal{P}]$ imply the newly deduced sentence $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}])$. Suppose that $\mathcal{F}[\mathcal{P}]$ and $\mathcal{G}[\mathcal{P}]$ are true; we would like to show that then $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}])$ is true. We show that one of the two disjuncts, $\mathcal{F}[\text{false}]$ or $\mathcal{G}[\text{true}]$, is true.

In the case in which the common subsentence \mathcal{P} is false, we know (by the *value* property, because \mathcal{P} and *false* have the same truth value and $\mathcal{F}[\mathcal{P}]$ is true) that the first of the disjuncts, $\mathcal{F}[\text{false}]$, is true.

Similarly, in the case in which the common subsentence \mathcal{P} is true, we know (by the *value* property again, because \mathcal{P} and *true* have the same truth value and $\mathcal{G}[\mathcal{P}]$ is true) that the second of the disjuncts, $\mathcal{G}[\text{true}]$, is true. \blacksquare

We have established the soundness of the ground version of the resolution rule when applied to ground sentences, which contain no variables. We require the sentences to be ground because the justification depends on the *value* property, which holds only for ground sentences. We can actually apply the ground version of the rule to sentences with variables; the soundness of such applications follows from the justification for the general version of the rule, which we present later.

We now discuss an important strategy for controlling the resolution rule.

THE POLARITY STRATEGY

Murray's [82] *polarity strategy* allows us to consider only those applications of the resolution rule under which at least one occurrence of \mathcal{P} is positive (or of no polarity) in $\mathcal{F}[\mathcal{P}]$ and at least one occurrence of \mathcal{P} is negative (or of no polarity) in $\mathcal{G}[\mathcal{P}]$. In other words, not all the subsentences that are replaced with *false* are negative and not all the subsentences that are replaced with *true* are positive. This strategy blocks many useless applications of the rule and rarely interferes with a reasonable step.

The intuitive rationale for the polarity strategy is that it is our goal to deduce the sentence *false*, which is more false than any other sentence. By replacing positive sentences with *false* and negative sentences with *true*, we are moving in the right direction, making the entire sentence more false.

Example

Suppose our deduced set contains the sentences

$$\mathcal{F}: \boxed{p(a)}^+ \text{ or } q(b)$$

and

$$\mathcal{G}: \text{if } \boxed{p(a)}^- \text{ then } q(b).$$

These sentences have occurrences of a common subsentence $p(a)$, of positive and negative polarity, respectively, as indicated by the annotation. By application of the resolution rule, we obtain the sentence

$$\begin{array}{c} \text{false or } q(b) \\ \text{or} \\ \text{if true then } q(b), \end{array}$$

which reduces to $q(b)$ under transformation.

Let us reverse the roles of our sentences.

$$\mathcal{F}: \text{if } \boxed{p(a)}^- \text{ then } q(b)$$

$$\mathcal{G}: \boxed{p(a)}^+ \text{ or } q(b).$$

The sentences still have occurrences of a common subsentence $p(a)$. However, it is in violation of the polarity strategy to apply the rule for the sentences in this order, because now the occurrence of $p(a)$ is negative in \mathcal{F} , i.e., it is not positive or of no polarity. Also, the polarity of $p(a)$ is positive in \mathcal{G} . If we insist on applying the resolution rule anyway, we obtain the sentence

$$\begin{array}{c} \text{if false then } q(b) \\ \text{or} \\ \text{true or } q(b), \end{array}$$

which reduces to *true* under transformation. Although it does no harm to add the sentence *true* to our deduced set, it is of no use in establishing the unsatisfiability of the set.

There are two other legal applications of the resolution rule to the same two sentences, obtained by taking the common subsentence to be $q(b)$ rather than $p(a)$. Both of these applications of the rule lead us to obtain the redundant sentence *true*, and both are in violation of the polarity strategy. \blacksquare

RESOLUTION RULE: GENERAL VERSION

The general version of the rule allows us to instantiate the variables of the given sentences as necessary to create common subsentences. It is expressed as follows:

Rule (resolution, general version)

For any sentences \mathcal{P} , $\tilde{\mathcal{P}}$, $\mathcal{F}[\mathcal{P}]$, and $\mathcal{G}[\tilde{\mathcal{P}}]$, where \mathcal{F} and \mathcal{G} are standardized apart, i.e., they have no variables in common, we have

$$\frac{\mathcal{F}[\mathcal{P}] \quad \mathcal{G}[\tilde{\mathcal{P}}]}{\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta[\text{true}]}$$

where θ is a most-general unifier of \mathcal{P} and $\tilde{\mathcal{P}}$.

More precisely,

- \mathcal{F} has one or more subsentences $\mathcal{P}, \mathcal{P}_1, \mathcal{P}_2, \dots$
- \mathcal{G} has one or more subsentences $\tilde{\mathcal{P}}, \tilde{\mathcal{P}}_1, \tilde{\mathcal{P}}_2, \dots$
- θ is a most general unifier of $\mathcal{P}, \mathcal{P}_1, \mathcal{P}_2, \dots$, and $\tilde{\mathcal{P}}, \tilde{\mathcal{P}}_1, \tilde{\mathcal{P}}_2, \dots$; hence

$$\mathcal{P}\theta = \mathcal{P}_1\theta = \mathcal{P}_2\theta = \dots = \tilde{\mathcal{P}}\theta = \tilde{\mathcal{P}}_1\theta = \tilde{\mathcal{P}}_2\theta = \dots$$

- The conclusion of the rule is obtained by replacing all occurrences of $\mathcal{P}\theta$ in $\mathcal{F}\theta$ with *false*, replacing all occurrences of $\tilde{\mathcal{P}}\theta$ (that is, $\mathcal{P}\theta$) in $\mathcal{G}\theta$ with *true*, and taking the disjunction of the results.

In other words, we apply the ground version of the rule to $\mathcal{F}\theta$ and $\mathcal{G}\theta$, taking $\mathcal{P}\theta$ as the common subsentence. \blacksquare

The rule requires that the sentences \mathcal{F} and \mathcal{G} be standardized apart, i.e., that they have no variables in common. This may be achieved by renaming the variables of the sentences as necessary. If both are the same sentence, we rename the variables of one copy of the sentence.

Let us show that the general version of the rule is sound.

Justification (resolution rule, general version):

The soundness of the general version of the rule follows from the soundness of its ground version. We show that the sentences \mathcal{F} and \mathcal{G} imply the sentence $(\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta[\text{true}])$.

We suppose that [under a given interpretation] the sentences \mathcal{F} and \mathcal{G} are true and show that $(\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta[\text{true}])$ is also true. It suffices (by the definition of truth for a nonground sentence) to show that any ground instance of $(\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta[\text{true}])$ is true.

Because \mathcal{F} and \mathcal{G} are true, we know (by the *instantiation lemma*) that $\mathcal{F}\theta$ and $\mathcal{G}\theta$ are true and hence (by the definition of truth for a nonground sentence) that every ground instance of $\mathcal{F}\theta$ and $\mathcal{G}\theta$ is true. But any ground instance of $(\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta[\text{true}])$ is the result of applying the ground version of the rule to the corresponding ground instance of $\mathcal{F}\theta$ and $\mathcal{G}\theta$; therefore it is also true. \blacksquare

The general version of the rule includes the ground version as a special case, in which the most-general unifier θ is the empty substitution $\{ \}$.

The following illustration of the general resolution rule is extracted from the derivation of a binary-search real-number square-root program.

Example

In the theory of the nonnegative real numbers, suppose our deduced set contains the sentence

$$\mathcal{F}: \text{not}(y^2 \leq a \text{ and not } \boxed{(y + \epsilon)^2 \leq a}^+),$$

where y is a variable and a and ϵ are constants. (The sentence is negated because it is deduced from the negation of the original theorem.)

We are about to apply the resolution rule to this sentence and itself. Therefore let us produce another copy of the sentence and standardize the two sentences apart; i.e., we rename the variable of the second sentence

$$\mathcal{G}: \text{not}(\boxed{\tilde{y}^2 \leq a} \text{ and not } ((\tilde{y} + \epsilon)^2 \leq a)).$$

The boxed subsentences

$$\mathcal{P}: (y + \epsilon)^2 \leq a$$

and

$$\tilde{\mathcal{P}}: \tilde{y}^2 \leq a$$

are unifiable, with most-general unifier

$$\theta: \{\tilde{y} \leftarrow y + \epsilon\}.$$

To apply the rule, we replace all occurrences of $\mathcal{P}\theta$ in $\mathcal{F}\theta$ with *false*, replace all occurrences of $\tilde{\mathcal{P}}\theta$ in $\mathcal{G}\theta$ with *true*, and take the disjunction of the results, obtaining

$$\begin{aligned} &\text{not}(y^2 \leq a \text{ and not false}) \\ &\text{or} \\ &\text{not}(true \text{ and not } (((y + \epsilon) + \epsilon)^2 \leq a)). \end{aligned}$$

This sentence reduces under transformation to

$$\text{not}(y^2 \leq a) \text{ or } ((y + \epsilon) + \epsilon)^2 \leq a.$$

The above application of the rule is in accordance with the polarity strategy, because the boxed subsentence \mathcal{P} is positive in \mathcal{F} and the boxed subsentence $\tilde{\mathcal{P}}$ is negative in \mathcal{G} . \blacksquare

The resolution rule presented here is an extension of the rule of Robinson [65] to the nonclausal case. Robinson's rule applies to clauses of the form

$$\begin{aligned} \mathcal{F}: & \mathcal{P} \text{ or } \mathcal{F}' \\ \mathcal{G}: & (\text{not } \tilde{\mathcal{P}}) \text{ or } \mathcal{G}', \end{aligned}$$

where \mathcal{P} and $\tilde{\mathcal{P}}$ are unifiable propositions, with most-general unifier θ , and \mathcal{F}' and \mathcal{G}' are themselves clauses. Robinson's rule deduces the new sentence

$$\mathcal{F}'\theta \text{ or } \mathcal{G}'\theta.$$

The resolution rule presented here deduces, from the same sentences \mathcal{F} and \mathcal{G} , the new sentence

$$\begin{aligned} &\text{false or } \mathcal{F}'\theta \\ &\text{or} \\ &(\text{not true}) \text{ or } \mathcal{G}'\theta. \end{aligned}$$

This sentence reduces under transformation to $(\mathcal{F}'\theta \text{ or } \mathcal{G}'\theta)$, the same sentence deduced by Robinson's version of the rule.

Nonclausal resolution was developed independently by Manna and Waldinger [80] and Murray [82]. The resolution and transformation rules together have been shown by Murray to provide a complete system for first-order logic. An implementation of a nonclausal resolution theorem prover by Stickel [82] employs a connection graph strategy.

5. THE RELATION REPLACEMENT RULE

We now begin to extend our nonclausal deduction system to give special treatment to a binary relation \prec . The two new rules of the extension allow us to build into the system instances of the *polarity replacement* proposition, just as the paramodulation and E-resolution rules allow us to build in instances of the substitutivity of equality.

Recall that, according to the *polarity replacement* proposition, for any sentence $P\langle x^+, y^- \rangle$ and binary relation \prec , the sentence

$$\begin{array}{l} \text{if } x \prec y \\ \text{then if } P\langle x^+, y^- \rangle \text{ then } P\langle y^+, x^- \rangle \end{array}$$

is valid.

If we could add this sentence to our deduced set for each relevant sentence $P\langle x^+, y^- \rangle$, we could achieve a considerable abbreviation of the proof, at the cost of a dramatic explosion of the search space. The extended system will behave as if the sentences were present, achieving the same abbreviation of the proof and, at the same time, collapsing rather than exploding the search space.

We begin with the relation replacement rule, which is our generalization of the paramodulation rule.

THE GROUND VERSION

With respect to a given relation \prec , the rule allows us to replace subexpression occurrences with larger or smaller expressions, depending on their polarity. The ground version of the rule which applies to sentences with no variables, is as follows:

Rule (relation replacement, ground version)

For any binary relation \prec , ground expressions s and t , and ground sentences $F[s \prec t]$ and $G\langle s^+, t^- \rangle$, we have

$$\frac{\begin{array}{c} F[s \prec t] \\ G\langle s^+, t^- \rangle \end{array}}{F[\text{false}] \text{ or } G\langle t^+, s^- \rangle}.$$

Here $G\langle t^+, s^- \rangle$ is obtained from $G\langle s^+, t^- \rangle$ by replacing certain positive occurrences of s with t and replacing certain negative occurrences of t with s , where polarity is taken in $G\langle s^+, t^- \rangle$ with respect to \prec . \blacksquare

In other words, if $F[s \prec t]$ and $G\langle s^+, t^- \rangle$ are sentences in our deduced set, we can add to the set the sentence $(F[\text{false}] \text{ or } G\langle t^+, s^- \rangle)$.

For a particular relation \prec , we shall refer to this rule as the \prec -replacement rule: thus, we have a $<$ -replacement rule, a \leq -replacement rule, and so forth. Although the rule allows us to replace occurrences in

$\mathcal{G}(s^+, t^-)$ of both expressions s and t at the same time, it is typically applied to replace occurrences of one or the other expression, but not both. Subsequent application of transformation rules, to remove occurrences of the truth symbols *true* and *false*, may be regarded as part of the relation replacement rule itself.

There is a polarity strategy for the relation replacement rule, which allows us to apply the rule only if some occurrence of $s \prec t$ is positive (or of no polarity) in $\mathcal{F}[s \prec t]$.

Naturally we may also require that some occurrence of s or t is actually replaced; otherwise, $\mathcal{G}(t^+, s^-)$ is identical to $\mathcal{G}(s^+, t^-)$, and the sentence we obtain is $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}(s^+, t^-))$; this is weaker than the sentence $\mathcal{G}(s^+, t^-)$, which was already in the deduced set.

In illustrating the rule we draw boxes around the matching occurrences of s and t .

Example

In the theory of the nonnegative integers, suppose our deduced set contains the sentences

$$\mathcal{F}: \begin{array}{l} \text{if } p(s) \\ \text{then } (\boxed{s} < t)^+ \end{array}$$

and

$$\mathcal{G}: s < \boxed{s^+}^2.$$

Note that the boxed occurrence of s in \mathcal{G} is positive with respect to the less-than relation $<$. Therefore we can apply the $<$ -replacement rule to replace the occurrence of s in \mathcal{G} with t , to deduce

$$\left[\begin{array}{l} \text{if } p(s) \\ \text{then false} \end{array} \right] \text{ or } s < t^2,$$

which reduces under transformation to

$$(\text{not } p(s)) \text{ or } s < t^2.$$

The above application of the rule is in accordance with the polarity strategy, because the occurrence of $s < t$ is positive in \mathcal{F} . Note that not every occurrence of s in \mathcal{G} was replaced in applying the rule.

In a system without the relation replacement rule, we could have deduced the same conclusion by applying the resolution rule in sequence to \mathcal{F} , \mathcal{G} , the *monotonicity* property

$$\begin{array}{l} \text{if } x < y \\ \text{then } x^2 < y^2, \end{array}$$

and the *transitivity* property

$$\begin{array}{l} \text{if } x < y \\ \text{then if } y < z \\ \text{then } x < z. \end{array}$$

The rule allows us to draw the conclusion even if the *monotonicity* and *transitivity* properties are not in our deduced set. \perp

The following illustration of the rule is extracted from the derivation of a program to find the maximum element of a list of numbers.

Example

In a theory of lists of numbers (integers, say), suppose our deduced set contains the sentences

$$\mathcal{F}: \begin{array}{l} \text{not} \left[\begin{array}{l} \text{if } g(m) = h \\ \text{then not } (m < \boxed{h})^+ \end{array} \right] \\ \text{or} \\ t = [] \end{array}$$

and

$$\mathcal{G}: \text{not} \left[\begin{array}{l} \text{if } g(h) \in t \\ \text{then } g(h) \leq \boxed{h}^- \end{array} \right]$$

Note that the boxed occurrence of h in \mathcal{G} is negative with respect to $<$. Therefore we can apply the $<$ -replacement rule to replace the occurrence of h in \mathcal{G} with m , to deduce

$$\begin{array}{l} \left[\begin{array}{l} \text{not} \left[\begin{array}{l} \text{if } g(m) = h \\ \text{then not false} \end{array} \right] \\ \text{or} \\ t = [] \end{array} \right] \\ \text{or} \\ \text{not} \left[\begin{array}{l} \text{if } g(h) \in t \\ \text{then } g(h) \leq m \end{array} \right] \end{array}$$

This sentence reduces under *true-false* transformation to

$$\begin{array}{l} t = [] \\ \text{or} \\ \text{not} \left[\begin{array}{l} \text{if } g(h) \in t \\ \text{then } g(h) \leq m \end{array} \right] \end{array}$$

The above application of the rule is in accordance with the polarity strategy, because the subsentence $m < h$ is positive in \mathcal{F} . \perp

Let us now establish the soundness of the rule.

Justification (relation replacement, ground version)

We show that the given sentences $\mathcal{F}[s \prec t]$ and $\mathcal{G}(s^+, t^-)$ imply the conclusion $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}(t^+, s^-))$. We distinguish between two cases and show that in each case one of the two disjuncts, $\mathcal{F}[\text{false}]$ or $\mathcal{G}(t^+, s^-)$, is true.

In the case in which the subsentence $s \prec t$ is false, we know (by the *value* property, because $s \prec t$ and *false* have the same truth value and $\mathcal{F}[s \prec t]$ is true) that the first of the disjuncts, $\mathcal{F}[\text{false}]$, is true.

In the case in which $s \prec t$ is true, we know (by the *polarity replacement* proposition, because $\mathcal{G}(s^+, t^-)$ is true) that the second of the disjuncts, $\mathcal{G}(t^+, s^-)$, is true. \perp

As with the resolution rule, we have established the soundness of the ground version of the relation replacement rule when applied to sentences with no variables. We will actually apply the ground version of the rule to sentences with variables. The above justification does not extend to this case, however, because the *value* property only holds for ground sentences. Such applications are an instance of the following general version of the rule.

THE GENERAL VERSION

We are now ready to give the general version of the rule, which applies to sentences with variables and allows us to instantiate the variables as necessary to create common subexpressions.

Rule (relation replacement, general version)

For any binary relation \prec , expressions s, t, \tilde{s} , and \tilde{t} , and sentences $\mathcal{F}[s \prec t]$ and $\mathcal{G}(\tilde{s}^+, \tilde{t}^-)$, where \mathcal{F} and \mathcal{G} are standardized apart, we have

$$\frac{\mathcal{F}[s \prec t] \quad \mathcal{G}(\tilde{s}^+, \tilde{t}^-)}{\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta(t\theta^+, s\theta^-)}$$

where θ is a simultaneous, most-general unifier of s, \tilde{s} and of t, \tilde{t} .

More precisely,

- \mathcal{F} has one or more subsentences $s \prec t, s_1 \prec t_1, s_2 \prec t_2, \dots$
- \mathcal{G} has one or more subexpressions $\tilde{s}, \tilde{s}_1, \tilde{s}_2, \dots$ and $\tilde{t}, \tilde{t}_1, \tilde{t}_2, \dots$
- θ is a simultaneous most-general unifier of $s, s_1, s_2, \dots, \tilde{s}, \tilde{s}_1, \tilde{s}_2, \dots$ and of $t, t_1, t_2, \dots, \tilde{t}, \tilde{t}_1, \tilde{t}_2, \dots$; hence

$$s\theta = s_1\theta = s_2\theta = \dots = \tilde{s}\theta = \tilde{s}_1\theta = \tilde{s}_2\theta = \dots$$

and

$$t\theta = t_1\theta = t_2\theta = \dots = \tilde{t}\theta = \tilde{t}_1\theta = \tilde{t}_2\theta = \dots$$

- The conclusion of the rule is obtained by replacing all occurrences of $(s \prec t)\theta$ in $\mathcal{F}\theta$ with *false*, replacing certain positive occurrences of $s\theta$ in $\mathcal{G}\theta$ with $t\theta$, replacing certain negative occurrences of $t\theta$ in $\mathcal{G}\theta$ with $s\theta$, and taking the disjunction of the two results. Here polarity is in $\mathcal{G}\theta$ with respect to \prec .

In other words, we apply the ground version of the rule to $\mathcal{F}\theta$ and $\mathcal{G}\theta$. \blacksquare

The justification of the general version of the rule, which we omit, is straightforward now that the soundness of the ground version has been established. The proof is analogous to the proof of the general version of the resolution rule. The polarity strategy for this rule allows us to assume that at least one occurrence of the subsentence $(s \prec t)\theta$ is positive or of no polarity in $\mathcal{F}\theta$.

Example

In the theory of sets, suppose our deduced set contains the sentences

$$\mathcal{F}: \begin{array}{l} \text{if } p(x) \\ \text{then } (\boxed{h(x, a)} \subset \boxed{b})^+ \text{ or } (\boxed{h(b, y)} \subset \boxed{x})^+ \end{array}$$

and

$$\mathcal{G}: (c \in \boxed{h(u, a)})^+ \sim v \text{ or } q(u, v),$$

where \sim is the set difference function.

Note that

- \mathcal{F} contains the [positive] subsentences $h(x, a) \subset b$ and $h(b, y) \subset x$.
- The boxed subterms $h(x, a)$, $h(b, y)$, and $h(u, a)$ and the boxed subterms b and x are simultaneously unifiable, with most-general unifier

$$\theta : \{x \leftarrow b, u \leftarrow b, y \leftarrow a\}.$$
- The boxed occurrence of $h(u, a)$ is positive in \mathcal{G} with respect to \subset .

Therefore we can apply the \subset -replacement rule, replacing all occurrences of $h(b, a) \subset b$ in $\mathcal{F}\theta$ with *false*, replacing the occurrence of $h(b, a)$ in $\mathcal{G}\theta$ with b , and taking the disjunction of the results, to obtain

$$\begin{array}{l} \left[\begin{array}{l} \text{if } p(b) \\ \text{then false or false} \end{array} \right] \\ \text{or} \\ (c \in b \sim v) \text{ or } q(b, v). \end{array}$$

This sentence reduces under transformation to

$$(not\ p(b)) \text{ or } (c \in b \sim v) \text{ or } q(b, v).$$

The above application of the rule is in accordance with the polarity strategy. \blacksquare

Use of the relation replacement rule allows a dramatic abbreviation of many proofs. For this reason and because the rule enables us to eliminate troublesome axioms from the deduced set, the search space is constricted. We have not established completeness results for the rule; judging from the corresponding theorem for paramodulation (Brand [75]), we expect such results to be difficult.

SPECIAL CASE: THE EQUALITY REPLACEMENT RULE

The most important instance of the relation replacement rule is obtained by taking the relation \prec to be the equality relation $=$. This special case of the rule, which allows us to replace equals with equals, is a nonclausal version of the paramodulation rule. It may be expressed as follows:

Rule (equality replacement)

For any terms s , t , \tilde{s} , and \tilde{t} , and sentences $\mathcal{F}[s = t]$ and $\mathcal{G}(\tilde{s}, \tilde{t})$, where \mathcal{F} and \mathcal{G} are standardized apart, we have

$$\frac{\begin{array}{c} \mathcal{F}[s = t] \\ \mathcal{G}(\tilde{s}, \tilde{t}) \end{array}}{\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta(t\theta, s\theta)}$$

where θ is a simultaneous, most-general unifier of s , \tilde{s} and of t , \tilde{t} . \blacksquare

The notation is analogous to that for the general relation-replacement rule. We do not need to restrict the polarity of the replaced subterms $s\theta$ and $t\theta$ in $\mathcal{G}\theta$, because any term has both polarities with respect to the equality relation. The polarity strategy is the same as before.

The following illustration of the equality replacement rule is extracted from the derivation of an integer quotient program.

Example

In the theory of the nonnegative integers, suppose our deduced set contains the sentences

$$\mathcal{F}: (\boxed{0 \cdot u} = 0)^+$$

and

$$\mathcal{G}: \text{not} (\boxed{z \cdot d} \leq n \text{ and } (z + 1) \cdot d > n).$$

(In the derivation, \mathcal{F} is an axiom and \mathcal{G} is deduced from the negation of the theorem.)

Note that

- \mathcal{F} contains the (positive) subsentence $0 \cdot u = 0$.
- The boxed subterms $0 \cdot u$ and $z \cdot d$ are unifiable, with most-general unifier

$$\theta: \{z \leftarrow 0, u \leftarrow d\}.$$

Therefore we can apply the $=$ -replacement rule, replacing all occurrences of $0 \cdot d = 0$ in $\mathcal{F}\theta$ with *false*, replacing the occurrence of $0 \cdot d$ in $\mathcal{G}\theta$ with 0, and taking the disjunction of the results, to deduce

$$\begin{array}{c} \text{false} \\ \text{or} \\ \text{not} (0 \leq n \text{ and } (0 + 1) \cdot d > n). \end{array}$$

This sentence reduces under *true-false* transformation to

$$\text{not} (0 \leq n \text{ and } (0 + 1) \cdot d > n). \quad \blacksquare$$

SPECIAL CASE: THE EQUIVALENCE REPLACEMENT RULE

Another important instance of the relation replacement rule is obtained by taking the relation \Leftarrow to be the equivalence connective \equiv . This is possible only because we regard connectives as relations over truth values. The rule is analogous to the equality replacement rule.

Rule (equivalence replacement rule)

For any sentences $S, \tau, \tilde{S}, \tilde{\tau}, \mathcal{F}[S \equiv \tau]$, and $\mathcal{G}(\tilde{S}, \tilde{\tau})$, where \mathcal{F} and \mathcal{G} are standardized apart, we have

$$\frac{\mathcal{F}[S \equiv \tau] \quad \mathcal{G}(\tilde{S}, \tilde{\tau})}{\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta(\tau\theta, S\theta)}$$

where θ is a simultaneous, most-general unifier of S, \tilde{S} and of $\tau, \tilde{\tau}$. \blacksquare

As in the equality replacement rule, we do not need to restrict the polarities of the replaced subsentences $S\theta$ and $\tau\theta$ in $\mathcal{G}\theta$, because any subsentence has both polarities with respect to the equivalence relation. The polarity strategy is the same as for the general relation-replacement rule.

The following illustration of the equivalence replacement rule (or \equiv -replacement rule) is drawn from the derivation of a program to find the maximum of a list of numbers (e.g., integers or reals).

Example

In the theory of lists of (say) integers, suppose our deduced set contains the sentences

$$\mathcal{F}: \begin{array}{l} \text{if not } (x = \{ \}) \\ \text{then } \boxed{u \in x} \equiv [u = h \text{ or } u \in t] \end{array}^+$$

and

$$\mathcal{G}: \text{not } \left[\begin{array}{l} z \in s \text{ and} \\ \text{if } \boxed{g(z) \in s} \\ \text{then } z \geq g(z) \end{array} \right]$$

(In the derivation, \mathcal{F} is an axiom and \mathcal{G} is deduced from the negation of the theorem.)

Note that the boxed subsentences $u \in x$ and $g(z) \in s$ are unifiable, with most-general unifier

$$\theta: \{u \leftarrow g(z), x \leftarrow s\}.$$

Therefore we can apply the \equiv -replacement rule, replacing the occurrence of $g(z) \in s$ in $\mathcal{G}\theta$ with

$$g(z) = h \text{ or } g(z) \in t,$$

to deduce

$$\begin{array}{l} \left[\begin{array}{l} \text{if not } (s = \{ \}) \\ \text{then false} \end{array} \right] \\ \text{or} \\ \text{not } \left[\begin{array}{l} z \in s \text{ and} \\ \text{if } [g(z) = h \text{ or } g(z) \in t] \\ \text{then } z \geq g(z) \end{array} \right] \end{array}$$

This sentence reduces under transformation to

$$\begin{array}{l} s = \{ \} \\ \text{or} \\ \text{not } \left[\begin{array}{l} z \in s \text{ and} \\ \text{if } [g(z) = h \text{ or } g(z) \in t] \\ \text{then } z \geq g(z) \end{array} \right] \end{array}$$

6. THE RELATION-MATCHING RULE

We are about to introduce not a rule in itself but an augmentation of the other rules. The resolution and relation replacement rules draw a conclusion when one subexpression in our proof unifies with another. The relation-matching augmentation allows these rules to apply even if the two expressions fail to unify, provided that certain conditions can be introduced into the conclusion. We begin by describing the augmentation of the resolution rule.

RESOLUTION WITH RELATION MATCHING: GROUND VERSION

This rule is our generalization of the E-resolution rule. The ground version of the rule is as follows:

Rule (resolution with relation matching, ground version)

For any binary relation \prec , ground expressions s and t , and ground sentences $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$, $\mathcal{F}[\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle]$, and $\mathcal{G}[\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle]$ we have

$$\frac{\mathcal{F}[\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle] \quad \mathcal{G}[\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle]}{\text{if } s \preceq t \text{ then } \mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}]}$$

Here

- $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ is an arbitrary sentence, called the *intermediate* sentence, which may have positive and negative occurrences of s and t ; polarity is taken with respect to \prec .
- The sentence \mathcal{F} may have several distinct subsentences $\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle$, each obtained from the intermediate sentence $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ by replacing certain of the positive occurrences of t with s and certain of the negative occurrences of s with t .
- Similarly, \mathcal{G} may have several distinct subsentences $\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle$, each obtained from the intermediate sentence by replacing certain of the positive occurrences of s with t and certain of the negative occurrences of t with s . \blacksquare

For a particular relation \prec , we shall refer to the above as the resolution rule with \prec -matching.

Note that if all the subsentences $\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle$ and $\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle$ were identical, we could apply the original resolution rule, obtaining the conclusion $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}])$. The augmented rule allows us to derive the same conclusion rule even if the subsentences \mathcal{P} do not match exactly, provided that the mismatches occur between terms s and t of restricted polarity and that the condition $s \preceq t$ is introduced.

The polarity strategy allows us to apply the rule only if an occurrence of one of the sentences $\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle$ is positive or of no polarity in \mathcal{F} and if an occurrence of one of the sentences $\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle$ is negative or of no polarity in \mathcal{G} .

Note that the intermediate sentence $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ does not necessarily appear in either of the sentences of the deduced set and that the rule does not stipulate how to find such a sentence. We shall discuss the choice of the intermediate sentence in the subsection **Selection of Application Parameters**.

Example

In the theory of lists, suppose that our deduced set includes the sentences

$$\mathcal{F}: p(\ell) \text{ or } \boxed{c \in (\text{tail}(\ell))^+}^+$$

and

$$\mathcal{G}: \text{if } \boxed{c \in \ell^+} \text{ then } q(\ell).$$

The two boxed subsentences are not identical. Let us take our intermediate sentence to be one of them, $\mathcal{P}: c \in \text{tail}(\ell)$. The subterm $s^+ : \text{tail}(\ell)$ is positive in $c \in \text{tail}(\ell)$ with respect to the proper-sublist relation \prec_{list} . The other boxed subsentence $c \in \ell$ can be obtained by replacing this subterm with $t^+ : \ell$. Therefore we can apply the resolution rule with \prec_{list} -matching to obtain

$$\begin{aligned} &\text{if } \text{tail}(\ell) \preceq_{\text{list}} \ell \\ &\text{then } p(\ell) \text{ or false} \\ &\quad \text{or} \\ &\text{if true then } q(\ell), \end{aligned}$$

which reduces under transformation to

$$\begin{array}{l} \text{if } \text{tail}(\ell) \leq_{\text{list}} \ell \\ \text{then } p(\ell) \text{ or } q(\ell). \quad \blacksquare \end{array}$$

We shall give some more complex examples of the application of the rule after we establish its soundness.

Justification (resolution with relation matching, ground version)

Note that (by the invertibility of partial replacement) the intermediate sentence $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ can be obtained from any of the subsentences $\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle$ of \mathcal{F} by replacing certain positive occurrences of s with t and certain negative occurrences of t with s , where polarity is taken in \mathcal{P} with respect to \prec . Therefore (by the *polarity replacement* proposition) each of the sentences

$$\begin{array}{l} \text{if } s \preceq t \\ \text{then if } \mathcal{P}\langle s^+, s^+, t^-, t^- \rangle \\ \quad \text{then } \mathcal{P}\langle s^+, t^+, s^-, t^- \rangle \end{array} \quad (\dagger)$$

is valid.

Also any of the subsentences $\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle$ of \mathcal{G} can be obtained from the intermediate sentence $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ by replacing certain positive occurrences of s with t and certain negative occurrences of t with s . Therefore (by the *polarity replacement* proposition again) each of the sentences

$$\begin{array}{l} \text{if } s \preceq t \\ \text{then if } \mathcal{P}\langle s^+, t^+, s^-, t^- \rangle \\ \quad \text{then } \mathcal{P}\langle t^+, t^+, s^-, s^- \rangle \end{array} \quad (\ddagger)$$

is valid.

Suppose that the sentences $\mathcal{F}[\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle]$ and $\mathcal{G}[\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle]$ are true and that $s \preceq t$. We would like to show that then $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}])$ is true. The proof distinguishes between two cases, depending on whether the intermediate sentence $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ is false or true. We show that in each case one of the two disjuncts, $\mathcal{F}[\text{false}]$ or $\mathcal{G}[\text{true}]$, is true.

Case: $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ is false

Then by our previous conclusion (\dagger) , because $s \preceq t$, we know each of the subsentences $\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle$ of \mathcal{F} is false. Because $\mathcal{F}[\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle]$ is true and because the subsentences $\mathcal{P}\langle s^+, s^+, t^-, t^- \rangle$ and false all have the same truth value, we know (by the *value* property) that the first disjunct, $\mathcal{F}[\text{false}]$, is true.

Case: $\mathcal{P}\langle s^+, t^+, s^-, t^- \rangle$ is true

Then by our previous conclusion (\ddagger) , because $s \preceq t$, we know each of the sentences $\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle$ is true. Because $\mathcal{G}[\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle]$ is true and because $\mathcal{P}\langle t^+, t^+, s^-, s^- \rangle$ and true have the same truth value, we know (by the *value* property again) that the second disjunct, $\mathcal{G}[\text{true}]$, is true. \blacksquare

The resolution rule with relation matching must be regulated with strict heuristic controls; if the controls are too permissive, any two subsentences may be matched.

The following example is a bit contrived but illustrates some of the power of the rule.

Example

In the theory of sets, suppose our deduced set includes the two sentences

$$\mathcal{F}: \begin{array}{c} \boxed{e \in ((s^+ \sim a) \cup (b \sim t^-) \cup (t^+ \sim c) \cup (d \sim t^-))}^+ \\ \text{or} \\ \boxed{e \in ((s^+ \sim a) \cup (b \sim s^-) \cup (s^+ \sim c) \cup (d \sim t^-))}^+ \end{array}$$

and

$$\mathcal{G}: \text{not} \left[\begin{array}{c} \boxed{e \in ((t^+ \sim a) \cup (b \sim s^-) \cup (t^+ \sim c) \cup (d \sim t^-))}^- \\ \text{and} \\ \boxed{e \in ((s^+ \sim a) \cup (b \sim s^-) \cup (t^+ \sim c) \cup (d \sim s^-))}^- \end{array} \right]$$

Let us take our intermediate sentence to be

$$\mathcal{P}: e \in ((s^+ \sim a) \cup (b \sim s^-) \cup (t^+ \sim c) \cup (d \sim t^-)).$$

The occurrences of s and t have been annotated with their polarities in \mathcal{P} with respect to the proper-subset relation \subset . Note that each of the boxed sentences in \mathcal{F} may be obtained from \mathcal{P} by replacing certain of the positive occurrences of t with s and certain of the negative occurrences of s with t . Also, each of the boxed subsentences of \mathcal{G} may be obtained from \mathcal{P} by replacing certain of the positive occurrences of s with t and certain of the negative occurrences of t with s . Therefore we can apply the resolution rule with \subset -matching to obtain

$$\begin{array}{l} \text{if } s \subseteq t \\ \text{then false or false} \\ \text{or} \\ \text{not (true and true),} \end{array}$$

which reduces under transformation to the sentence

$$\text{not } (s \subseteq t). \quad \lrcorner$$

Note that this conclusion, obtained by a single application of the rule, is not immediately evident to the human reader.

SPECIAL CASE: RESOLUTION WITH EQUALITY MATCHING

In the case in which the relation \prec is taken to be the equality relation $=$, the resolution rule with relation matching reduces to a nonclausal variant of the E-resolution rule. It may be expressed (in the ground version) as follows:

Rule (resolution with equality matching)

For any terms s and t and sentences $\mathcal{P}(s, t, s, t)$, $\mathcal{F}[\mathcal{P}(s, s, t, t)]$, and $\mathcal{G}[\mathcal{P}(t, t, s, s)]$, we have

$$\frac{\mathcal{F}[\mathcal{P}(s, s, t, t)] \quad \mathcal{G}[\mathcal{P}(t, t, s, s)]}{\begin{array}{l} \text{if } s = t \\ \text{then } \mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}]. \quad \lrcorner \end{array}}$$

Here $\mathcal{P}\langle s, s, t, t \rangle$ and $\mathcal{P}\langle t, t, s, s \rangle$ are obtained from $\mathcal{P}\langle s, t, s, t \rangle$ by replacing certain occurrences of s with t and certain occurrences of t with s . In other words, all the subsentences $\mathcal{P}\langle s, s, t, t \rangle$ and $\mathcal{P}\langle t, t, s, s \rangle$ are identical except that one may have occurrences of s where another has occurrences of t . We do not need to restrict the polarities, because every subterm of a sentence is both positive and negative with respect to the equality relation.

MULTIPLE MISMATCHED SUBSENTENCES

The resolution rule with relation matching can be extended to allow several corresponding pairs of subexpressions $s_1, t_1, s_2, t_2, \dots$ and s_n, t_n rather than a single pair s, t , and several binary relations \prec_1, \prec_2, \dots , and \prec_n rather than a single binary relation \prec . To write the extended rule succinctly, we abbreviate s_1, s_2, \dots, s_n as $\hat{s}, t_1, t_2, \dots, t_n$ as $\hat{t}, \prec_1, \prec_2, \dots$, and \prec_n as $\hat{\prec}$, and

$$s_1 \prec_1 t_1 \text{ and } s_2 \prec_2 t_2 \text{ and } \dots \text{ and } s_n \prec_n t_n \text{ as } \hat{s} \hat{\prec} \hat{t}.$$

Then for any binary relations $\hat{\prec}$, expressions \hat{s} and \hat{t} , and sentences $\mathcal{F}[\mathcal{P}\langle \hat{s}^+, \hat{t}^+, \hat{s}^-, \hat{t}^- \rangle]$, $\mathcal{G}[\mathcal{P}\langle \hat{t}^+, \hat{t}^+, \hat{s}^-, \hat{s}^- \rangle]$, we have

$$\frac{\mathcal{F}[\mathcal{P}\langle \hat{s}^+, \hat{s}^+, \hat{t}^-, \hat{t}^- \rangle] \quad \mathcal{G}[\mathcal{P}\langle \hat{t}^+, \hat{t}^+, \hat{s}^-, \hat{s}^- \rangle]}{\text{if } \hat{s} \hat{\prec} \hat{t} \text{ then } \mathcal{F}[\text{false}] \text{ or } \mathcal{G}[\text{true}]}.$$

The extended rule is easily justified, given the soundness of the original rule.

RESOLUTION WITH RELATION MATCHING: GENERAL VERSION

The general version of the rule allows us to instantiate the variables of the given sentences as necessary and then to apply the ground version. The precise statement, which we omit, is analogous to the precise statement of the general version of the resolution rule. We illustrate the application of the general rule with an example.

Example

Suppose our deduced set contains the sentences

$$\mathcal{F}: \quad \text{if } q(u) \text{ then } \boxed{p(u^+, u^+)}^+$$

and

$$\mathcal{G}: \quad \text{not } \boxed{p(\ell^+, f(\ell)^+)}^-.$$

Here the annotations of the subterms within the boxed subsentences indicate their polarity in these subsentences with respect to a binary relation \prec .

The substitution $\theta: \{u \leftarrow \ell\}$ fails to unify the boxed subsentences of \mathcal{F} and \mathcal{G} ; the results of applying θ to these subsentences are the sentences $p(\ell^+, \ell^+)$ and $p(\ell^+, f(\ell)^+)$, respectively. Note that the mismatched occurrences of ℓ and $f(\ell)$ are positive in these sentences with respect to \prec .

To apply the ground version of the rule to $\mathcal{F}\theta$ and $\mathcal{G}\theta$, let us take the intermediate sentence to be $p(\ell^+, \ell^+)$. We obtain

$$\begin{array}{l} \text{if } \ell \preceq f(\ell) \\ \text{then } \left[\begin{array}{l} \text{if } q(\ell) \\ \text{then false} \end{array} \right] \text{ or (not true),} \end{array}$$

which reduces under *true-false* transformation to

$$\begin{array}{l} \text{if } \ell \preceq f(\ell) \\ \text{then not } q(\ell). \quad \lrcorner \end{array}$$

SELECTION OF APPLICATION PARAMETERS

For each application of the resolution rule with relation matching, we must select the *application parameters*, i.e., the substitution θ , the intermediate sentence \mathcal{P} , and the subexpressions s and t . In fact, a satisfactory choice of application parameters is not straightforward: it depends on what other sentences are in the deductive set. Some considerations influencing the decision are illustrated in the next few sections.

Choice of Substitution

The substitution θ and the intermediate sentence \mathcal{P} for applying the rule are not necessarily unique.

In the example above, consider again the boxed subsentences $p(u^+, u^+)$ and $p(\ell^+, f(\ell)^+)$ of \mathcal{F} and \mathcal{G} . Instead of the substitution $\theta : \{u \leftarrow \ell\}$, consider the substitution $\theta' : \{u \leftarrow f(\ell)\}$. This substitution also fails to unify the boxed subsentences; the results of applying θ' to the boxed subsentences are the sentences $p(f(\ell)^+, f(\ell)^+)$ and $p(\ell^+, f(\ell)^+)$, respectively. Note that the mismatched occurrences of $f(\ell)$ and ℓ are positive in these sentences with respect to \prec .

To apply the ground version of the rule to $\mathcal{F}\theta'$ and $\mathcal{G}\theta'$, let us take the intermediate sentence to be $p(f(\ell)^+, f(\ell)^+)$. We obtain

$$\begin{array}{l} \text{if } f(\ell) \preceq \ell \\ \text{then } \left[\begin{array}{l} \text{if } q(\ell) \\ \text{then false} \end{array} \right] \text{ or (not true),} \end{array}$$

which reduces under *true-false* transformation to

$$\begin{array}{l} \text{if } f(\ell) \preceq \ell \\ \text{then not } q(\ell). \end{array}$$

This is not equivalent to the sentence we obtained by applying the rule with the substitution θ ,

$$\begin{array}{l} \text{if } \ell \preceq f(\ell) \\ \text{then not } q(\ell). \end{array}$$

In other words, we must consider both ways of applying the rule.

To Unify or Not to Unify

In previous examples, we have applied the resolution rule with relation matching only when it is illegal to apply the ordinary resolution rule because the matched subsentences fail to unify. In some cases, however, we must use relation matching to obtain a refutation even though the matched subsentences do unify and the resolution rule could be applied.

For example, suppose our deduced set consists of the sentences

1. $\boxed{p(x^+)} \text{ or } q(x^+)$
2. $\text{not } \boxed{p(a^+)}^-$
3. $\text{not } \boxed{q(b^+)}^-$
4. $c \preceq a$
5. $c \preceq b,$

where x is positive in the boxed subsentence $p(x)$ and in the subsentence $q(x)$ with respect to the relation \preceq , as indicated by its annotation.

It is legal to apply the ordinary resolution rule to the first two sentences, taking the unifier to be $\{x \leftarrow a\}$, to deduce (after transformation)

$$q(a).$$

However, this sentence is of no use in a refutation.

If instead we apply the resolution rule with \preceq -matching to the same boxed subsentences, taking the unifier to be the empty substitution $\{ \}$, we obtain (after transformation)

6. $\text{if } x \preceq a \text{ then } \boxed{q(x^+)}^+.$

We can then apply the resolution rule to sentences 6 and 3, taking the unifier to be the empty substitution $\{ \}$, to obtain (after transformation)

7. $\text{if } x \preceq b \text{ then not } (x \preceq a).$

We finally obtain a refutation by applying the resolution rule to this sentence and the last two sentences in turn; the unifier is $\{x \leftarrow c\}$.

In applying the ordinary resolution rule, we committed x to be a ; this turned out to be a mistake. In applying the resolution rule with \preceq -matching instead, we left x free to be any element such that $x \preceq a$; in particular, we could then take x to be c .

Choice of Mismatched Subexpressions

In the examples of resolution with relation matching we have seen, we have always taken the mismatched subexpressions s and t to be as small as possible. Sometimes this choice costs us a proof.

For instance, suppose our deduced set consists of the sentences

1. $\boxed{p(f(a))}^+$
2. $\text{not } \boxed{p(f(b))}^-$
3. $f(a) = f(b).$

If we apply the resolution rule with equality matching to the first two sentences, taking s to be a and t to be b , we obtain

$$\begin{array}{l} \text{if } a = b \\ \text{then false or not true,} \end{array}$$

which reduces under transformation to

$$\text{not } (a = b).$$

This sentence is of no use in a refutation.

On the other hand, if instead we apply the same rule taking s to be $f(a)$ and t to be $f(b)$, we obtain

if $f(a) = f(b)$
then false or not true,

which reduces under transformation to

not $(f(a) = f(b))$.

A refutation can be obtained immediately by applying the resolution rule to the third sentence and this one.

In the preceding examples, we have seen that in applying the resolution rule with relation matching, the choice of appropriate application parameters, i.e., the substitution θ , the intermediate sentence \mathcal{P} , and the mismatched subexpressions s and t , are not unique and depend on the other sentences in the deduced set. Digricoli [83] provides an algorithm to generate all legal sets of application parameters. This algorithm is phrased in terms of his variant of the E-resolution rule but extends readily to the general, nonclausal case. Digricoli also suggests a heuristic *viability criterion* for selecting a single appropriate set of application parameters; this criterion appears to extend to the general case as well.

REPLACEMENT WITH RELATION MATCHING: GROUND VERSION

We have shown how to augment the resolution rule to apply even if the matched subsentences are not entirely unified by the substitution. We now introduce an analogous augmentation of the relation replacement rule.

Rule (replacement with relation matching, ground version)

For any binary relations \prec_1 and \prec_2 , ground expressions $s, t, u\langle s^+, t^+, s^-, t^- \rangle$, and $v\langle s^+, t^+, s^-, t^- \rangle$, and ground sentences

$$\mathcal{F}[u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle]$$

and

$$\mathcal{G}\langle u\langle t^+, t^+, s^-, s^- \rangle^+, v\langle t^+, t^+, s^-, s^- \rangle^- \rangle,$$

we have

$$\frac{\mathcal{F}[u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle] \quad \mathcal{G}\langle u\langle t^+, t^+, s^-, s^- \rangle^+, v\langle t^+, t^+, s^-, s^- \rangle^- \rangle}{\text{if } s \preceq_2 t \text{ then } \mathcal{F}[\text{false}] \text{ or } \mathcal{G}\langle v\langle t^+, t^+, s^-, s^- \rangle^+, u\langle t^+, t^+, s^-, s^- \rangle^- \rangle}$$

Here

- The expressions $u\langle s^+, t^+, s^-, t^- \rangle$ and $v\langle s^+, t^+, s^-, t^- \rangle$ are arbitrary expressions. The sentence $u\langle s^+, t^+, s^-, t^- \rangle \prec_1 v\langle s^+, t^+, s^-, t^- \rangle$ is called the *intermediate sentence*.
- The subsentences $u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle$ of \mathcal{F} are obtained from the intermediate sentence by replacing certain positive occurrences of t with s and certain negative occurrences of s with t , where polarity is taken in the intermediate sentence with respect to \prec_2 .
- The subexpressions $u\langle t^+, t^+, s^-, s^- \rangle$ and $v\langle t^+, t^+, s^-, s^- \rangle$ of \mathcal{G} are obtained from $u\langle s^+, t^+, s^-, t^- \rangle$ and $v\langle s^+, t^+, s^-, t^- \rangle$, respectively, by replacing certain occurrences of s

with t and certain occurrences of t with s , where again polarity is taken in the intermediate sentence with respect to \prec_2 .

- The subsentence $\mathcal{G}(v(t^+, t^+, s^-, s^-)^+, u(t^+, t^+, s^-, s^-)^-)$ of the conclusion is obtained from $\mathcal{G}(u(t^+, t^+, s^-, s^-)^+, v(t^+, t^+, s^-, s^-)^-)$ by replacing certain positive occurrences of $u(t^+, t^+, s^-, s^-)$ with $v(t^+, t^+, s^-, s^-)$ and certain negative occurrences of $v(t^+, t^+, s^-, s^-)$ with $u(t^+, t^+, s^-, s^-)$, where the polarity of u and v is taken in \mathcal{G} with respect to \prec_1 . \blacksquare

For particular binary relations \prec_1 and \prec_2 , we shall call this the \prec_1 -replacement rule with \prec_2 -matching. Note that if $u(t^+, t^+, s^-, s^-)$ and $v(t^+, t^+, s^-, s^-)$ were identical to $u(s^+, s^+, t^-, t^-)$ and $v(s^+, s^+, t^-, t^-)$, respectively, we could apply the original \prec_1 -replacement rule without \prec_2 -matching, obtaining the conclusion

$$\mathcal{F}[\text{false}] \text{ or } \mathcal{G}(v(t^+, t^+, s^-, s^-)^+, u(t^+, t^+, s^-, s^-)^-).$$

The augmented rule allows us to derive the same conclusion, even if the subexpressions do not match exactly, provided that the mismatches occur between subexpressions s and t of restricted polarity with respect to \prec_2 and that the condition $s \preceq_2 t$ is added.

Example

In a theory that includes the lists and the integers, suppose our deduced set contains the sentences

$$\mathcal{F}: (\boxed{\text{length}(m^-)}^- \leq a) \text{ or } p(m)$$

and

$$\mathcal{G}: \text{if } q(\ell) \text{ then } (\boxed{\text{length}(\ell^-)}^+ > b),$$

where ℓ and m are lists and a and b are integers.

The two boxed subexpressions are not identical, so we cannot apply the original \leq -replacement rule. To apply the augmented rule, let us take our intermediate sentence to be $\text{length}(\ell) \leq a$. With respect to the proper sublist relation \prec_{list} , the subterm $s^- : \ell$ is negative in the intermediate sentence $u \prec_1 v : \text{length}(\ell) \leq a$. From this sentence we can obtain the subsentence $\text{length}(m) \leq a$ of \mathcal{F} by replacing the negative occurrence of ℓ with $t^- : m$. Therefore, by the \leq -replacement rule with \prec_{list} -matching, we deduce

$$\begin{aligned} &\text{if } \ell \preceq_{list} m \\ &\text{then false or } p(m) \\ &\quad \text{or} \\ &\text{if } q(\ell) \text{ then } a > b. \end{aligned}$$

Here the subsentence $a > b$ of the conclusion is obtained from the subsentence $\text{length}(\ell) > b$ of \mathcal{G} by replacing a positive occurrence of $u^+ : \text{length}(\ell)$ with $v^+ : b$, where polarity is taken in \mathcal{G} with respect to the weak less-than relation \leq . The conclusion reduces under transformation to

$$\begin{aligned} &\text{if } \ell \preceq_{list} m \\ &\text{then } p(m) \text{ or} \\ &\quad \text{if } q(\ell) \text{ then } a > b. \end{aligned} \quad \blacksquare$$

Now let us establish the soundness of the rule.

Justification (replacement with relation matching, ground version)

Note that (by the invertibility of partial replacements), the intermediate sentence $u(s^+, t^+, s^-, t^-) \prec_1 v(s^+, t^+, s^-, t^-)$ can be obtained from any of the subsentences $u(s^+, s^+, t^-, t^-) \prec_1 v(s^+, s^+, t^-, t^-)$ of

\mathcal{F} by replacing certain positive occurrences of s with t and certain negative occurrences of t with s , where polarity is taken in the subsentences with respect to \prec_2 . Therefore (by the *polarity replacement* proposition), each of the sentences

$$(†) \quad \begin{array}{l} \text{if } s \preceq_2 t \\ \text{then if } u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle \\ \quad \text{then } u\langle s^+, t^+, s^-, t^- \rangle \prec_1 v\langle s^+, t^+, s^-, t^- \rangle \end{array}$$

is valid.

Also any of the sentences $u\langle t^+, t^+, s^-, s^- \rangle \prec_1 v\langle t^+, t^+, s^-, s^- \rangle$ can be obtained from the intermediate sentence $u\langle s^+, t^+, s^-, t^- \rangle \prec_1 v\langle s^+, t^+, s^-, t^- \rangle$ by replacing certain positive occurrences of s with t and certain negative occurrences of t with s , where polarity is taken in the intermediate sentence with respect to \prec_2 . Therefore (by the *polarity replacement* proposition again) each of the sentences

$$(‡) \quad \begin{array}{l} \text{if } s \preceq_2 t \\ \text{then if } u\langle s^+, t^+, s^-, t^- \rangle \prec_1 v\langle s^+, t^+, s^-, t^- \rangle \\ \quad \text{then } u\langle t^+, t^+, s^-, s^- \rangle \prec_1 v\langle t^+, t^+, s^-, s^- \rangle \end{array}$$

is valid.

Furthermore the subsentence $\mathcal{G}\langle v\langle t^+, t^+, s^-, s^- \rangle^+, u\langle t^+, t^+, s^-, s^- \rangle^- \rangle$ of the conclusion can be obtained from the given sentence $\mathcal{G}\langle u\langle t^+, t^+, s^-, s^- \rangle^+, v\langle t^+, t^+, s^-, s^- \rangle^- \rangle$ of the deduced set by replacing certain positive occurrences of $u\langle t^+, t^+, s^-, s^- \rangle$ with $v\langle t^+, t^+, s^-, s^- \rangle$ and certain negative occurrences of $v\langle t^+, t^+, s^-, s^- \rangle$ with $u\langle t^+, t^+, s^-, s^- \rangle$, where polarity is taken in \mathcal{G} with respect to \prec_1 . Therefore (by the *polarity replacement* proposition once again) each of the sentences

$$(††) \quad \begin{array}{l} \text{if } u\langle t^+, t^+, s^-, s^- \rangle \prec_1 v\langle t^+, t^+, s^-, s^- \rangle \\ \text{then if } \mathcal{G}\langle u\langle t^+, t^+, s^-, s^- \rangle^+, v\langle t^+, t^+, s^-, s^- \rangle^- \rangle \\ \quad \text{then } \mathcal{G}\langle v\langle t^+, t^+, s^-, s^- \rangle^+, u\langle t^+, t^+, s^-, s^- \rangle^- \rangle \end{array}$$

is valid.

Suppose that the ground sentences

$$\mathcal{F}[u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle] \quad \text{and} \quad \mathcal{G}\langle u\langle t^+, t^+, s^-, s^- \rangle^+, v\langle t^+, t^+, s^-, s^- \rangle^- \rangle$$

are true and that $s \preceq_2 t$. We would like to show that then

$$\mathcal{F}[\text{false}] \text{ or } \mathcal{G}\langle v\langle t^+, t^+, s^-, s^- \rangle^+, u\langle t^+, t^+, s^-, s^- \rangle^- \rangle$$

is true. The proof distinguishes between two cases, depending on whether the intermediate sentence is false or true. We show that in each case one of the two disjuncts, $\mathcal{F}[\text{false}]$ or $\mathcal{G}\langle v\langle t^+, t^+, s^-, s^- \rangle^+, u\langle t^+, t^+, s^-, s^- \rangle^- \rangle$, is true.

Case: $u\langle s^+, t^+, s^-, t^- \rangle \prec_1 v\langle s^+, t^+, s^-, t^- \rangle$ is false

Then by our previous conclusion (†), because $s \preceq_2 t$, we know each of the subsentences $u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle$ of \mathcal{F} is false. Because $\mathcal{F}[u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle]$ is true and because the sentences $u\langle s^+, s^+, t^-, t^- \rangle \prec_1 v\langle s^+, s^+, t^-, t^- \rangle$ and *false* all have the same truth value, we know (by the *value* property) that the first disjunct, $\mathcal{F}[\text{false}]$, is true.

Case: $u\langle s^+, t^+, s^-, t^- \rangle \prec_1 v\langle s^+, t^+, s^-, t^- \rangle$ is true

Then by our previous conclusion (‡), because $s \preceq_2 t$, we know each of the sentences $u\langle t^+, t^+, s^-, s^- \rangle \prec_1 v\langle t^+, t^+, s^-, s^- \rangle$ is true. Therefore by several applications of our previous conclusion (††), because

$$\mathcal{G}\langle u\langle t^+, t^+, s^-, s^- \rangle^+, v\langle t^+, t^+, s^-, s^- \rangle^- \rangle$$

is true, we know that the second disjunct,

$$\mathcal{G}\langle v\langle t^+, t^+, s^-, s^- \rangle^+, u\langle t^+, t^+, s^-, s^- \rangle^- \rangle,$$

is true.

In each case, we have shown that the desired conclusion is true. \square

REPLACEMENT WITH RELATION MATCHING: GENERAL VERSION

The general version of the rule allows us to instantiate the variables of the given sentences as necessary and then to apply the ground version. We omit the precise statement, which is analogous to the general version of the relation replacement rule, but we illustrate the general version with an example extracted from the derivation of a program to sort a list of numbers.

Example

In a theory of lists of (say) integers, suppose our deduced set contains the sentences

$$\mathcal{F} : \boxed{\text{perm}(x_1 \sqcup \langle u \rangle \sqcup x_2, y_1 \sqcup \langle u \rangle \sqcup y_2)} \equiv \text{perm}(x_1 \sqcup x_2, y_1 \sqcup y_2)}^+$$

and

$$\mathcal{G} : \text{not}(\text{ordered}(z) \text{ and } \boxed{\text{perm}(\ell^+, z)}).$$

Here the term $x_1 \sqcup x_2$ is the result of appending the lists x_1 and x_2 , and the term $\langle u \rangle$ is the list whose sole element is u . Also, $\text{perm}(\ell, z)$ holds if the list ℓ is a permutation of the list z , and $\text{ordered}(z)$ holds if the elements of z are in (weakly) increasing order. In the derivation, \mathcal{F} is one of the axioms for the permutation relation, which states that two lists are permutations if they are still permutations after dropping a common element, and \mathcal{G} is the negation of the theorem, which states the existence of an ordered list that is a permutation of a given list.

The results of applying the substitution

$$\theta : \{z \leftarrow y_1 \sqcup \langle u \rangle \sqcup y_2\}$$

to the boxed subsentences are

$$\text{perm}((x_1 \sqcup \langle u \rangle \sqcup x_2)^+, y_1 \sqcup \langle u \rangle \sqcup y_2)$$

and

$$\text{perm}(\ell^+, y_1 \sqcup \langle u \rangle \sqcup y_2).$$

The mismatched subterms

$$x_1 \sqcup \langle u \rangle \sqcup x_2 \quad \text{and} \quad \ell$$

are positive in their respective subsentences with respect to the perm relation. (Because this relation is symmetric, they also happen to be negative.) The boxed subsentence $\text{perm}(\ell, z)$ is positive in \mathcal{G} with respect to the equivalence relation \equiv . (It also happens to be negative.) Therefore, by the \equiv -replacement rule with perm -matching, we may deduce the sentence

$$\begin{aligned} &\text{if } \text{perm}(x_1 \sqcup \langle u \rangle \sqcup x_2, \ell) \\ &\text{then false} \\ &\quad \text{or} \\ &\text{not}(\text{ordered}(y_1 \sqcup \langle u \rangle \sqcup y_2) \text{ and } \text{perm}(x_1 \sqcup x_2, y_1 \sqcup y_2)) \end{aligned}$$

which reduces under transformation to

$$\begin{array}{l} \text{if } \text{perm}(x_1 \sqcap (\langle u \rangle \sqcap x_2), \ell) \\ \text{then } \text{not } (\text{ordered}(y_1 \sqcap (\langle u \rangle \sqcap y_2)) \text{ and } \text{perm}(x_1 \sqcap x_2, y_1 \sqcap y_2)). \quad \lrcorner \end{array}$$

RELATION MATCHING VERSUS RELATION REPLACEMENT

The relation matching and relation replacement rules play complementary roles, and one might expect that a single deductive system would employ one or the other rule but not both. After all, in clausal equality systems, paramodulation and a variant of E-resolution have each been shown to be complete (Anderson [70], Digricoli [83], and Brand [75]) without including the other. Moreover, by incorporating both rules, we admit a troublesome redundancy: The same conclusion can be derived in several ways.

On the other hand, it often turns out that a proof that seems unmotivated or tricky using only one of the rules seems more straightforward using a combination of both. For instance, in an example of a previous section, we applied the resolution rule with relation matching to the sentences

$$\mathcal{F}: \begin{array}{l} \text{if } q(u) \\ \text{then } \boxed{p(u^+, u^+)}^+ \end{array}$$

and

$$\mathcal{G}: \text{not } \boxed{p(\ell^+, f(\ell)^+)}^-$$

taking the substitution to be

$$\theta: \{u \leftarrow \ell\},$$

to obtain after transformation

$$\begin{array}{l} \text{if } \boxed{\ell \prec f(\ell)} \\ \text{then } \text{not } q(\ell). \end{array}$$

If our deduced set also contains the sentence

$$\boxed{v \prec f(v)},$$

we can further deduce (by resolution) the sentence

$$\text{not } q(\ell).$$

Now suppose our deductive system includes the relation replacement rule but not the relation-matching rule. Then to deduce the same conclusion $\text{not } q(\ell)$, we would have to apply the relation replacement rule to the sentences

$$\boxed{v}^+ \prec f(v)$$

and

$$\mathcal{G}: \text{not } p(\boxed{\ell^+}, f(\ell))$$

to obtain (after transformation)

$$\text{not } \boxed{p(f(\ell), f(\ell))}$$

We could then obtain the same conclusion ($\text{not } q(\ell)$) by resolution applied to this sentence and the sentence

$$\mathcal{F}: \begin{array}{ll} \text{if } q(u) \\ \text{then } \boxed{p(u, u)}. \end{array}$$

Although both sequences of inference lead to the same conclusion, the earlier proof seems better motivated: Each step is based on matching subexpressions that already possess a high degree of syntactic similarity. In contrast, the above proof seems rather gratuitous: The application of the relation replacement rule is based on matching the variable v with the constant ℓ . There is no reason to perform this step except as a preparation for the subsequent resolution step.

Examples can also be exhibited for which a proof employing the replacement rule is well-motivated but the corresponding proof using the matching rule appears strained. For instance, in the theory of integers, use of the $=$ -replacement rule and the axiom $u + (-u) = 0$ allows us to simplify a subterm of form $t + (-t)$ to 0. If we are only permitted to use the relation-matching rule, we must leave the subterm intact, and hope that we attempt to match it against a corresponding subterm 0 later in the proof.

We expect that by including both rules together in a system we shall be able to apply more restrictive strategies to each of them. Consequently, we shall obtain a smaller search space than if we had included either of the rules separately.

7. STRENGTHENING

The relation replacement rule of Section 5 does not always allow us to draw the strongest possible conclusion. In this section we establish a stronger form of the polarity replacement lemma and use it to develop a stronger relation-replacement rule.

We motivate the strengthening of the rule with an example. In the theory of the integers, suppose our deduced set contains the sentences

$$\mathcal{F}: \boxed{s} < t$$

and

$$\mathcal{G}: a \leq \boxed{s}^+ + 2.$$

Because the occurrence of s in \mathcal{G} is positive with respect to the less-than relation $<$, the $<$ -replacement rule allows us to replace s with t and deduce that (after transformation)

$$a \leq t^+ + 2.$$

From these two sentences, however, we should be able to deduce the stronger result

$$a < t + 2.$$

Similarly, from the sentence $s < t$ and $\text{not } (a - s > b)$, we should be able to deduce $\text{not } (a - t \geq b)$ rather than merely $\text{not } (a - t > b)$.

Unfortunately, the rule as we have presented it does not yield these more useful conclusions; the strengthened relation-replacement rule will. But first, we must introduce some preliminary notions.

THE STRENGTHENED POLARITY-REPLACEMENT LEMMA

The strengthened rule depends on the following basic result:

Lemma (strengthened polarity replacement)

Consider arbitrary expressions $e\langle x, y \rangle$ and $e'\langle x, y \rangle$ and binary relations \prec_1 and \prec_2 . The sentence

$$\begin{array}{l} \text{if } x \prec_1 y \\ \text{then if } e\langle x, y \rangle \prec_2 e'\langle x, y \rangle \\ \text{then } e\langle y, x \rangle \prec_2 e'\langle y, x \rangle \end{array}$$

is valid provided that the replaced occurrences of x and y satisfy the following *strengthening conditions* [in $e\langle x, y \rangle$ and $e'\langle x, y \rangle$ with respect to \prec_1 and \prec_2]:

- *transitivity condition*

The relation \prec_2 , the irreflexive restriction of \prec_2 , is transitive.

- *top condition*

The replaced occurrences of x and y are respectively positive and negative in $e\langle x, y \rangle \prec_2 e'\langle x, y \rangle$ with respect to \prec_1 .

- *left-right condition*

One of the following two disjuncts holds:

The replaced occurrences of x and y in $e\langle x, y \rangle$ are respectively negative and positive in $e\langle x, y \rangle$ with respect to \prec_1 and \prec_2 (and some replacement is made in $e\langle x, y \rangle$)

(*left disjunct*)

or

the replaced occurrences of x and y in $e'\langle x, y \rangle$ are respectively positive and negative in $e'\langle x, y \rangle$ with respect to \prec_1 and \prec_2 (and some replacement is made in $e'\langle x, y \rangle$).

(*right disjunct*) \blacksquare

Before proving this proposition, let us illustrate it with an example.

Example (strengthened polarity-replacement lemma)

In a theory that includes the sets and the nonnegative integers, take \prec_1 to be the proper-subset relation \subset over the sets and \prec_2 to be the weak less-than relation \leq over the nonnegative integers. Then \prec_2 is the strict less-than relation $<$.

Consider the sentence

$$m \cdot \text{card}(y) \leq n + \text{card}(x),$$

where x and y are sets, m and n are nonnegative integers, and $\text{card}(x)$ is the cardinality of the set x . According to the lemma, the sentence

$$\begin{array}{l} \text{if } x \subset y \\ \text{then if } m \cdot \text{card}(y) \leq n + \text{card}(x) \\ \text{then } m \cdot \text{card}(x) < n + \text{card}(y) \end{array}$$

is valid, because the replaced occurrences of x and y satisfy the strengthening conditions in $m \cdot \text{card}(y)$ and $n + \text{card}(x)$ with respect to \subset and \leq . In particular,

- The relation $<$ is transitive; hence the *transitivity condition* is satisfied.

- The replaced occurrences of x and y are respectively positive and negative in $m \cdot \text{card}(y) \leq n + \text{card}(x)$ with respect to \subset ; hence the *top* condition is satisfied.
- Although the replaced occurrence of y is not positive in $m \cdot \text{card}(y)$ with respect to \subset and $<$ (after all, m could be 0), the replaced occurrence of x is positive in $n + \text{card}(x)$ with respect to \subset and $<$. Hence, though the *left* disjunct of the *left-right* condition is not satisfied, the *right* disjunct is. \blacksquare

We are now ready to establish the lemma.

Proof (strengthened polarity-replacement lemma)

Suppose that

$$x \prec_1 y \quad \text{and} \quad e\langle x, y \rangle \prec_2 e'\langle x, y \rangle,$$

and that the strengthening conditions are satisfied.

We would like to show that then

$$e\langle y, x \rangle \prec_2 e'\langle y, x \rangle.$$

The *left-right* condition was stated as a disjunction of two possibilities; we treat each possibility separately.

Case (left disjunct): The replaced occurrences of x and y in $e\langle x, y \rangle$ are respectively negative and positive in $e\langle x, y \rangle$ with respect to \prec_1 and \prec_2 (and some replacement is made in $e\langle x, y \rangle$).

In this case (by the *transitive polarity-replacement* lemma, because $x \prec_1 y$), we have

$$e\langle y, x \rangle \prec_2 e\langle x, y \rangle.$$

Also (by the *polarity replacement* proposition and our supposition that $x \prec_1 y$ and $e\langle x, y \rangle \prec_2 e'\langle x, y \rangle$) we have

$$e\langle x, y \rangle \prec_2 e'\langle y, x \rangle.$$

(Here we have only performed the replacements on the right-hand side; by the *top* condition, we know the replaced occurrences of x and y are respectively positive and negative in $e\langle x, y \rangle \prec_2 e'\langle x, y \rangle$ with respect to \prec_1 .) It follows that

$$e\langle x, y \rangle \prec_2 e'\langle y, x \rangle \quad \text{or} \quad e\langle x, y \rangle = e'\langle y, x \rangle.$$

Because $e\langle y, x \rangle \prec_2 e\langle x, y \rangle$, we thus have (either by the transitivity of \prec_2 or the substitutivity of equality) that

$$e\langle y, x \rangle \prec_2 e'\langle y, x \rangle,$$

as we wanted to show.

Case (right disjunct): The replaced occurrences of x and y in $e'\langle x, y \rangle$ are respectively positive and negative in $e'\langle x, y \rangle$ with respect to \prec_1 and \prec_2 (and some replacement is made in $e\langle x, y \rangle$).

The proof in this case is entirely symmetric to the proof in the previous case. \blacksquare

THE STRENGTHENED POLARITY-REPLACEMENT PROPOSITION

The strengthened rule is expressed in terms of the following notational device:

Definition (strengthen accordingly)

Suppose \prec is a binary relation, s and t are expressions (either both sentences or both terms), and \mathcal{G} is a sentence.

If we write \mathcal{G} as $\mathcal{G}(s^+, t^-)$, then $\mathcal{G}(t^+, s^-)^\dagger$ denotes the sentence obtained by replacing certain positive occurrences of s with t , replacing certain negative occurrences of t with s (where polarity is taken with respect to \prec), and *strengthening accordingly* as follows:

- Whenever a replacement is made in a positive subsentence of form $e(s, t) \tilde{\prec} e'(s, t)$, where the replaced occurrences of s and t satisfy the strengthening conditions in $e(s, t)$ and $e'(s, t)$ with respect to \prec and $\tilde{\prec}$, replace the occurrence of the symbol $\tilde{\prec}$ with $\tilde{\prec}$, the irreflexive restriction of $\tilde{\prec}$.
- Whenever a replacement is made in a negative subsentence of form $e(s, t) \tilde{\prec} e'(s, t)$, where the replaced occurrences of s and t satisfy the strengthening conditions in $e(s, t)$ and $e'(s, t)$ with respect to \prec and $\tilde{\prec}$, replace the occurrence of the symbol $\tilde{\prec}$ with $\tilde{\prec}$. (Here $\tilde{\prec}$ and $\tilde{\prec}$ are the negation, and the reflexive closure, respectively, of $\tilde{\prec}$.) \blacksquare

These conditions may appear mysterious at this point, but they are precisely what we need to establish the following result, which tightens up the polarity replacement proposition:

Proposition (strengthened polarity replacement)

For any binary relation \prec and sentence $\mathcal{P}(x^+, y^-)$, the sentence

$$\begin{array}{l} \text{if } x \prec y \\ \text{then if } \mathcal{P}(x^+, y^-) \\ \text{then } \mathcal{P}(y^+, x^-)^\dagger \end{array}$$

is valid. \blacksquare

We illustrate the proposition with two examples.

Example

In the theory of the positive integers (excluding 0), take \prec to be the proper-divides relation \prec_{div} and take our sentence to be

$$\mathcal{P}(x^+, y^-) : a \leq (x+1)^2 \text{ or } q(x).$$

Then according to the proposition, the sentence

$$\begin{array}{l} \text{if } x \prec_{div} y \\ \text{then if } a \leq (x+1)^2 \text{ or } q(x) \\ \text{then } a < (y+1)^2 \text{ or } q(x) \end{array}$$

is valid. Note that the symbol \leq has been replaced by its irreflexive restriction $<$ as a result of the strengthening. This is because

- The subsentence $a \leq (x+1)^2$ is positive in $\mathcal{P}(x^+, y^-)$.
- The replaced occurrence of x in $a \leq (x+1)^2$ satisfies the strengthening conditions in a and $(x+1)^2$ with respect to \prec_{div} and \leq . In particular
 - The relation $<$ is transitive; hence the *transitivity* condition is satisfied.

- The replaced occurrence of x is positive in $a \leq (x+1)^2$ with respect to \prec_{div} ; hence the *top* condition is satisfied.
- The replaced occurrence of x is positive in $(x+1)^2$ with respect to \prec_{div} and $<$; hence the *right* disjunct of the *left-right* condition is satisfied. ─

Example

In a theory that includes the lists and the nonnegative integers, take \prec to be the tail relation \prec_{tail} over the lists and take our sentence to be

$$\mathcal{P}\langle x^+, y^- \rangle : \text{ if } length(x \sqcup \ell) < length(y) + m \text{ then } q(x, y),$$

where x , y , and ℓ are lists, m is a nonnegative integer, and $length(\ell)$ is the number of elements in the list ℓ . Then according to the proposition, the sentence

$$\begin{aligned} &\text{if } x \prec_{tail} y \\ &\text{then if } length(x \sqcup \ell) < length(y) + m \text{ then } q(x, y) \\ &\quad \text{then if } length(y \sqcup \ell) \leq length(x) + m \text{ then } q(x, y) \end{aligned}$$

is valid. Note that here the symbol $<$ has been replaced by \leq as a result of the strengthening. This is because

- The subsentence $length(x \sqcup \ell) < length(y) + m$ is negative in $\mathcal{P}\langle x^+, y^- \rangle$.
- The replaced occurrences of x and y satisfy the strengthening conditions in $length(x \sqcup \ell)$ and $length(y) + m$ with respect to \prec_{tail} and \prec , that is \geq . In particular
 - The relation $>$, the irreflexive restriction of \geq , is transitive; hence the *transitivity* condition is satisfied.
 - The replaced occurrences of x and y are positive and negative, respectively, in the sentence $length(x \sqcup \ell) \geq length(y) + m$ with respect to \prec_{tail} ; hence the *top* condition is satisfied.
 - The replaced occurrence of x is negative in $length(x \sqcup \ell)$ with respect to \prec_{tail} and $>$; hence the *left* disjunct of the *left-right* condition is satisfied. (As it turns out, the replaced occurrence of y is also negative in $length(y) + m$ with respect to \prec_{tail} and $>$; hence the *right* disjunct is also satisfied.) ─

Let us now prove the proposition.

Proof (strengthened polarity-replacement proposition)

We suppose that

$$x \prec y \quad \text{and} \quad \mathcal{P}\langle x^+, y^- \rangle,$$

and show that then

$$\mathcal{P}\langle y^+, x^- \rangle^\dagger.$$

The sentence $\mathcal{P}\langle y^+, x^- \rangle^\dagger$ is obtained from $\mathcal{P}\langle x^+, y^- \rangle$ by replacing certain subexpressions with others. We show that each of these replacements makes the sentence “truer,” in the sense that it produces a sentence implied by the original.

We consider separately three kinds of replacement:

- Replacing a positive subsentence of form $e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle$ with $e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle$, where the replaced occurrences of x and y satisfy the strengthening conditions in $e\langle x, y \rangle$ and $e'\langle x, y \rangle$ with respect to \prec and $\tilde{\sim}$.

In this case, because $x \prec y$, we have (by the *strengthened polarity-replacement lemma*) that

$$\begin{aligned} &\text{if } e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle \\ &\text{then } e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle. \end{aligned}$$

Therefore, because the replaced occurrence of $e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle$ is positive in $\mathcal{P}\langle x^+, y^- \rangle$, we know (by the original *polarity-replacement* proposition) that replacing it with the “truer” subsentence $e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle$ makes the entire sentence truer.

- Replacing a negative subsentence of form $e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle$, with $e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle$, where the replaced occurrences of x and y satisfy the strengthening conditions in $e\langle x, y \rangle$ and $e'\langle x, y \rangle$ with respect to \prec and $\tilde{\sim}$ (the negation of $\tilde{\sim}$).

In this case, because $x \prec y$, we have (by the *strengthened polarity-replacement lemma*, recalling that $\tilde{\sim}$ is the irreflexive restriction of $\tilde{\sim}$)

$$\begin{aligned} &\text{if } e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle \\ &\text{then } e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle \end{aligned}$$

or, equivalently (taking the contrapositive),

$$\begin{aligned} &\text{if } e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle \\ &\text{then } e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle. \end{aligned}$$

Therefore, because the replaced occurrence of $e\langle x, y \rangle \tilde{\sim} e'\langle x, y \rangle$ is negative in $\mathcal{P}\langle x^+, y^- \rangle$, we know (by the original *polarity-replacement* proposition) that replacing it with the “falsier” sentence $e\langle y, x \rangle \tilde{\sim} e'\langle y, x \rangle$ will make the entire sentence falsier.

- Replacing a positive occurrence of x with y or a negative occurrence of y with x , where polarity is with respect to \prec and where the replaced occurrence is not within the scope of any strengthened relation $\tilde{\sim}$.

In this case, the replacement makes the sentence “truer,” by the original *polarity-replacement* proposition. \blacksquare

THE GROUND VERSION

We can now express the stronger version of the relation replacement rule. The ground version of the rule is as follows:

Rule (strengthened relation replacement, ground version)

For any binary relation \prec , ground expressions s and t , and ground sentences $\mathcal{F}[s \prec t]$ and $\mathcal{G}\langle s^+, t^- \rangle$, we have

$$\frac{\mathcal{F}[s \prec t] \quad \mathcal{G}\langle s^+, t^- \rangle}{\mathcal{F}[\text{false}] \text{ or } \mathcal{G}\langle t^+, s^- \rangle^\dagger}$$

Here $\mathcal{G}(t^+, s^-)^\dagger$ is the result of replacing certain positive occurrences of s with t , replacing certain negative occurrences of t with s , and strengthening accordingly, where polarity is taken in $\mathcal{G}(s^+, t^-)$ with respect to \prec . We assume that at least one replacement is made. \blacksquare

Let us illustrate the ground version of the rule with two examples.

Example

In the theory of the positive integers (excluding 0), suppose our deduced set contains the sentences

$$\mathcal{F}: \text{ if } p(s) \text{ then } \boxed{s} \prec_{div} t$$

and

$$\mathcal{G}: a \leq (\boxed{s}^+ + 1)^2 \text{ or } q(s),$$

where \prec_{div} is the proper divides relation. Then we can apply the strengthened \prec_{div} -replacement rule to replace the boxed occurrence of s in \mathcal{G} with t and to strengthen accordingly, obtaining

$$\begin{array}{l} \text{if } p(s) \text{ then false} \\ \text{or} \\ a < (t + 1)^2 \text{ or } q(s). \end{array}$$

This sentence reduces under transformation to

$$(not\ p(s)) \text{ or } a < (t + 1)^2 \text{ or } q(s).$$

The relation symbol \leq was replaced by its irreflexive restriction $<$ because $a \leq (s + 1)^2$ is positive and because s and t satisfy the strengthening conditions in a and $(s + 1)^2$ with respect to \prec_{div} and \leq , as we have seen in a previous example.

Example

In a theory that includes the sets and the nonnegative integers, suppose our deduced set contains the sentences

$$\mathcal{F}: p(s, t) \text{ or } \boxed{s} \subset \boxed{t}$$

and

$$\mathcal{G}: not\ (q(s, t) \text{ and } m \cdot card(\boxed{s}^+) < n + card(\boxed{t}^-)),$$

where s and t are sets, m and n are nonnegative integers, and $card(s)$ is the cardinality of the set s . Then we can apply the strengthened \subset -replacement rule to replace the boxed occurrences of s with t and t with s and to strengthen accordingly, obtaining

$$\begin{array}{l} p(s, t) \text{ or false} \\ \text{or} \\ not\ (q(s, t) \text{ and } m \cdot card(t) \leq n + card(s)), \end{array}$$

that is (after transformation),

$$\begin{array}{l} p(s, t) \text{ or} \\ not\ (q(s, t) \text{ and } m \cdot card(t) \leq n + card(s)). \end{array}$$

The relation symbol $<$ has been replaced by its reflexive closure \leq because $m \cdot \text{card}(s) < n + \text{card}(t)$ is negative and because s and t satisfy the strengthening conditions in $m \cdot \text{card}(s)$ and $n + \text{card}(t)$ with respect to \subset and \nless , that is, \geq . In particular,

- The irreflexive restriction $>$ of \geq is transitive; hence the *transitivity* condition is satisfied.
- The replaced occurrences of s and t are respectively positive and negative in $m \cdot \text{card}(s) \leq n + \text{card}(t)$ with respect to \subset and \geq ; hence the *top* condition is satisfied.
- The replaced occurrence of t is negative in $n + \text{card}(t)$ with respect to \subset and $>$; hence the *right disjunct* of the *left-right* condition is satisfied. \blacksquare

Let us now establish the soundness of the rule.

Justification (relation replacement rule, ground version)

The proof resembles the justification of the original relation-replacement rule.

We suppose that the given sentences $\mathcal{F}[s \prec t]$ and $\mathcal{G}(s^+, t^-)$ are true and show that the newly deduced sentence $(\mathcal{F}[\text{false}] \text{ or } \mathcal{G}(t^+, s^-)^\dagger)$ is also true. We distinguish between two cases and show that in each case one of the two disjuncts, $\mathcal{F}[\text{false}]$ or $\mathcal{G}(t^+, s^-)^\dagger$, is true.

In the case in which the subsentence $s \prec t$ is false, we know (by the value property, because $s \prec t$ and *false* have the same truth value and $\mathcal{F}[s \prec t]$ is true) that the first of the disjuncts, $\mathcal{F}[\text{false}]$, is true.

In the case in which $s \prec t$ is true, we know (by the *strengthened polarity-replacement* proposition, because $\mathcal{G}(s^+, t^-)$ is true) that the second of the disjuncts, $\mathcal{G}(t^+, s^-)^\dagger$, is true. \blacksquare

THE GENERAL VERSION

The general version of the rule allows us to instantiate the variables of the sentences as necessary to create common subexpressions.

Rule (strengthened relation replacement, general version)

For any binary relation \prec , expressions s , t , \tilde{s} , and \tilde{t} , and sentences $\mathcal{F}[s \prec t]$ and $\mathcal{G}(\tilde{s}^+, \tilde{t}^-)$, where \mathcal{F} and \mathcal{G} are standardized apart, we have

$$\frac{\mathcal{F}[s \prec t] \quad \mathcal{G}(\tilde{s}^+, \tilde{t}^-)}{\mathcal{F}\theta[\text{false}] \text{ or } \mathcal{G}\theta(t\theta^+, s\theta^-)^\dagger},$$

where θ is a simultaneous, most-general unifier of s , \tilde{s} and of t , \tilde{t} . \blacksquare

As usual, to apply the general version of the rule to sentences \mathcal{F} and \mathcal{G} , we apply its ground version to $\mathcal{F}\theta$ and $\mathcal{G}\theta$. The justification, which is straightforward, is omitted. As before, the polarity strategy for the rule allows us to assume that a least one occurrence of the subsentence $(s \prec t)\theta$ is positive or of no polarity in $\mathcal{F}\theta$.

8. EXTENSIONS

The concepts in this paper are being extended in several directions. We briefly indicate several of these here.

EXPLICIT QUANTIFIERS

The system we have described deals with sentences that have had their quantifiers removed by skolemization. It is impossible, however, to remove quantifiers that occur within the scope of an equivalence (\equiv) connective or in the *if*-clause of a conditional (*if-then-else*) connective without first paraphrasing the connective in terms of others. If several of these connectives are nested, the paraphrased sentence becomes alarmingly complex.

In an earlier work (Manna and Waldinger [82]), we extend the deductive system to sentences that may have some of their quantifiers intact. In many cases, we can complete the proof without removing all the quantifiers. If these quantifiers are in equivalences or *if*-clauses, we need not paraphrase the offending connectives. Thus, we not only retain the form of the original sentence, but also can use the equivalences we retain in applying the equivalence replacement rule.

POLARITY WITH RESPECT TO AN EXPRESSION

We have used the notion of polarity with respect to a relation. Because a function is a special case of a relation, we can define polarity with respect to a function accordingly. Rather than restricting ourselves to the functions denoted by the function symbols in our deduced set, we prefer to consider the functions corresponding to particular expressions in the set.

Roughly speaking, suppose $e[s]$ is a ground term; then $e[s]$ corresponds to a binary relation $\prec_{e[s]}$ defined by the sentence

$$x \prec_{e[s]} y \equiv e[x] = y.$$

We may define polarity with respect to $\prec_{e[s]}$ just as we would with respect to any binary relation.

For example, in the theory of the integers, the relation $\prec_{e[s]}$ corresponding to the term $e[s] : s + 1$ is defined by the sentence

$$x \prec_{e[s]} y \equiv x + 1 = y.$$

(In fact, this relation turns out to be the predecessor relation \prec_{pred} we have seen earlier.) The relation $natnum(x)$, which holds if x is a nonnegative integer (natural number), is positive over its argument with respect to $\prec_{e[s]}$, for we have

$$\begin{array}{l} \text{if } x \prec_{e[s]} y \\ \text{then if } natnum(x) \\ \text{then } natnum(y). \end{array}$$

We can then establish an expression replacement rule analogous to our relation replacement rule; i.e., in the ground version:

For any expressions s and $e[s]$ and ground sentence $\mathcal{G}\langle s^+, e[s]^- \rangle$, we have

$$\frac{\mathcal{G}\langle s^+, e[s]^- \rangle}{\mathcal{G}\langle e[s]^+, s^- \rangle^\dagger}$$

Here $\mathcal{G}\langle e[s]^+, s^- \rangle^\dagger$ is obtained from $\mathcal{G}\langle s^+, e[s]^- \rangle$ by replacing certain positive occurrences of s with $e[s]$, replacing certain negative occurrences of $e[s]$ with s , and strengthening accordingly, where polarity is taken in $\mathcal{G}\langle s^+, e[s]^- \rangle$ with respect to $\prec_{e[s]}$.

For example, in the theory of the integers, if our deduced set contains the sentence

$$\mathcal{G}: \text{ not } [\text{natnum}((s+1)^-)]$$

we may deduce the sentence

$$\text{ not } [\text{natnum}(s)],$$

because the occurrence of $s+1$ is negative in \mathcal{G} with respect to the relation corresponding to the expression $s+1$.

We can also define expression-matching rules analogous to our relation-matching rule.

For example, in the theory of lists, suppose our deduced set contains the sentences

$$\mathcal{F}: \boxed{a \in s^+}$$

and

$$\mathcal{G}: \text{ not } (\boxed{a \in (b \circ s)^+}).$$

Here the term $b \circ s$ is the result of inserting the element b before the first element of the list s . By the resolution rule with expression matching, whose precise statement we omit, we may deduce (after transformation), the contradiction *false*, because s is positive in the boxed sentence $a \in s$ with respect to the relation corresponding to $b \circ s$.

CONDITIONAL POLARITY

Sometimes it is convenient to extend the notion of polarity to depend on the truth of certain conditions. For example, in the theory of integers (including negative integers) with respect to the relation \leq , the occurrence of s in the sentence

$$a < b \cdot s$$

might be regarded as positive if b is nonnegative and negative if b is nonpositive. (If b is 0, the occurrence might have both polarities). We could then adapt the relation replacement and relation matching rules to use this conditional polarity, imposing the appropriate conditions on whatever conclusion they draw.

More precisely, we define the notion of conditional polarity so that if x and y are respectively positive and negative in $\mathcal{P}\langle x^+, y^- \rangle$ with respect to the binary relation \prec subject to the condition $\mathcal{X}[x, y, \mathcal{Q}]$, then the sentence

$$\mathcal{X} \left[x, y, \begin{array}{l} \text{if } x \prec y \\ \text{then if } \mathcal{P}\langle x^+, y^- \rangle \\ \text{then } \mathcal{P}\langle y^+, x^- \rangle^\dagger \end{array} \right]$$

is valid. Here \mathcal{Q} denotes an arbitrary sentence; the indicated polarities of the replaced occurrences of x and y are subject to the condition $\mathcal{X}[x, y, \mathcal{Q}]$.

For example, according to this notion of conditional polarity, in the theory of the integers, the occurrence of x in the sentence

$$a \leq b + x^2$$

is positive with respect to the relation $<$ subject to the condition

$$\mathcal{N}[x, y, Q]: \begin{array}{l} \text{if } x \geq 0 \\ \text{then } Q. \end{array}$$

Consequently, we have that the sentence

$$\begin{array}{l} \text{if } x \geq 0 \\ \text{then if } x < y \\ \quad \text{then if } a \leq b + x^2 \\ \quad \quad \text{then } a < b + y^2 \end{array}$$

is valid. The relation \leq was replaced by $<$ as the result of strengthening.

In terms of this notion, we can introduce conditional versions of the relation replacement rule and relation-matching rules. In particular, we have the conditional relation-replacement rule, i.e., in the ground version:

For any binary relation \prec , ground expressions s and t , and ground sentences $\mathcal{F}[s \prec t]$ and $\mathcal{G}(s^+, t^-)$, we have

$$\frac{\mathcal{F}[s \prec t] \quad \mathcal{G}(s^+, t^-)}{\mathcal{N}[s, t, \text{false}] \text{ or } \mathcal{F}[\text{false}] \text{ or } \mathcal{G}(t^+, s^-)^{\uparrow}}.$$

Here the indicated polarities of the replaced occurrences of s and t are subject to the condition $\mathcal{N}[s, t, Q]$.

For example, in the theory of the integers, suppose our deduced set contains the sentences

$$\mathcal{F}: \begin{array}{l} \text{if } r(s, t) \\ \text{then } s < t \end{array}$$

and

$$\mathcal{G}: a < b \cdot s.$$

Note that the occurrence of s in \mathcal{G} is positive with respect to the relation $<$ subject to the condition

$$\begin{array}{l} \text{if } b \geq 0 \\ \text{then } Q. \end{array}$$

Therefore, according to the conditional $<$ -replacement rule, we may deduce

$$\left[\begin{array}{l} \text{if } b \geq 0 \\ \text{then false} \end{array} \right] \text{ or } \left[\begin{array}{l} \text{if } r(s, t) \\ \text{then false} \end{array} \right] \text{ or } a < b \cdot t,$$

which reduces under transformation to

$$(\text{not } (b \geq 0)) \text{ or } (\text{not } (r(s, t))) \text{ or } a < b \cdot t.$$

The conditional relation-matching rules are analogous. Of course these rules can be extended to apply to conditional polarity with respect to an expression rather than a relation.

PLANNING AND THE FRAME PROBLEM

Theorem-proving techniques have often been applied to problems in automatic planning. One approach to this application has been the formulation of a *situational logic*, a theory in which states of the world are

regarded as domain elements, denoted by terms. Typically, an action in a plan is represented as a function mapping states into other states. The effects of an action can be described by axioms.

For example, the primary effect of putting one block on top of another is expressed by an axiom such as

$$\begin{array}{l} \text{if } \text{clear}(x, w) \text{ and } \text{clear}(y, w) \\ \text{then } \text{on}(x, y, \text{puton}(x, y, w)). \end{array}$$

In other words, if block x is put on block y in a state w , then x will indeed be on y in the resulting state $\text{puton}(x, y, w)$. The antecedent expresses the preconditions that x and y be clear before x can be put on y ; in other words, no block can be on x or on y . (The conventional blocks-world hand can move only one block at a time.)

In a situational logic, a problem may be expressed as a theorem to be proved. For example, the problem of achieving the condition that block a is on block b and block b is on block c might be phrased as the theorem

$$(\exists z)[\text{on}(a, b, z) \text{ and } \text{on}(b, c, z)].$$

The *frame problem*, which occurs when planning problems are approached in this way, is connected with the requirement that we need to express not only what conditions are altered by a given action, but also what conditions are unchanged. For example, in addition to the primary effect of putting one block on top of another, we must state explicitly that this action has no effect on other relations, such as color; otherwise, we shall have no way of deducing that the color of a block after the action is the same as its color before. Therefore, we must include in our deduced set the *frame axiom*

$$\begin{array}{l} \text{if } \text{clear}(x, w) \text{ and } \text{clear}(y, w) \\ \text{then if } \text{color}(z, u, w) \\ \text{then } \text{color}(z, u, \text{puton}(x, y, w)). \end{array}$$

In other words, if the action of putting block x on top of block y is legal and if block z is of color u in state w , then z will also be of color u in the resulting state $\text{puton}(x, y, w)$. If our deduced set contains the sentence

$$\text{not } (\text{color}(c, \text{red}, \text{puton}(a, b, s))),$$

we can then apply the resolution rule to the frame axiom and this sentence to deduce (after transformation)

$$(\text{not } (\text{clear}(a, s))) \text{ or } (\text{not } (\text{clear}(b, s))) \text{ or } (\text{not } (\text{color}(c, \text{red}, s))).$$

We need a separate frame axiom not only for the color of blocks, but also their size, shape, surface texture, and any other attributes we wish to discuss in our theory. Adding all the frame axioms to our deduced set aggravates the search problem, because the axioms have many consequences irrelevant to the problem at hand.

By use of the conditional expression rules, we can drop all the frame axioms from our deduced set. For example, to paraphrase the above axiom we can declare that the relation $\text{color}(z, u, w)$ is positive with respect to the relation corresponding to the expression $e[w] : \text{puton}(x, y, w)$ subject to the condition

$$\mathcal{N}[w, w', \mathcal{Q}] : \begin{array}{l} \text{if } \text{clear}(x, w) \text{ and } \text{clear}(y, w) \\ \text{then } \mathcal{Q}. \end{array}$$

If our deduced set again contains the sentence

$$\text{not } (\text{color}(c, \text{red}, \text{puton}(a, b, s)^-)),$$

we can then apply the conditional expression-replacement rule to deduce

$$(\text{not } (\text{clear}(a, s))) \text{ or } (\text{not } (\text{clear}(b, s))) \text{ or } (\text{not } (\text{color}(c, \text{red}, s)))$$

as before, without requiring the frame axiom. Of course, the information that certain actions and relations are independent must still be expressed, but this can be done by polarity declarations rather than by axioms.

9. DISCUSSION

The theorem-proving system we have presented has been motivated by our work in program synthesis, and the best examples we have of its use are in this domain. We have used the system to write detailed derivations for programs over the integers and real numbers, the lists, the sets, and other structures. These derivations are concise and easy to follow: they reflect intuitive derivations of the same programs. A paper by Traugott [85] describes the application of this system to the derivation of several sorting programs. A paper by Manna and Waldinger [85] describes the derivation of several binary-search programs. Our earlier informal derivation of the unification algorithm (Manna and Waldinger [81]) can be expressed formally in this system.

An interactive implementation of the basic nonclausal theorem-proving system was completed by Malachi and has been extended by Bronstein to include some of the relation rules. An entirely automatic implementation is being contemplated. The relation rules will also be valuable for proving purely mathematical theorems. For this purpose they may be incorporated into clausal as well as nonclausal theorem-proving systems.

Theorem provers have exhibited superhuman abilities in limited subject domains, but seem least competent in areas in which human intuition is best developed. One reason for this is that an axiomatic formalization obscures the simplicity of the subject area; facts that a person would consider too obvious to require saying in an intuitive argument must be stated explicitly and dealt with in the corresponding formal proof. A person who is easily able to conduct the argument informally may well be unable to understand the formal proof, let alone to produce it.

Our work in special relations is part of a continuing effort to make formal theorem proving resemble intuitive reasoning. In the kind of system we envision, proofs are shorter, the search space is compressed, and heuristics based on human intuition become applicable.

ACKNOWLEDGEMENTS

The authors would like to thank Martin Abadi, Alex Bronstein, Tomas Feder, Eric Muller, Neil Murray, David Plaisted, Mark Stickel, Jon Traugott, and Frank Yellin for their suggestions and careful reading. Jon Traugott suggested extending the notion of polarity from one relation to two, making the rules more powerful and the exposition simpler; he also proposed the extended notions of polarity with respect to an expression and conditional polarity. The manuscript was prepared by Evelyn Eldridge-Diaz with the \TeX typesetting system.

REFERENCES

- Anderson [70]
R. Anderson, Completeness results for E-resolution, *AFIPS Spring Joint Computer Conference*, 1970, pp. 652-656.
- Boyer and Moore [79]
R. S. Boyer and J S. Moore, *A Computational Logic*, Academic Press, New York, N.Y., 1979.
- Brand [75]
D. Brand, Proving theorems with the modification method, *SIAM Journal of Computing*, Vol. 4, No. 2, 1975, pp. 412-430.
- Chang and Lee [73]

C. L. Chang and R. C. Lee, *Symbolic Logic and Mechanical Theorem Proving*, Academic Press, New York, N.Y., 1973.

Digricoli [83]

V. Digricoli, *Resolution By Unification and Equality*, Ph.D. thesis, New York University, New York, N.Y., 1983.

Kowalski [79]

R. Kowalski, *Logic for Problem Solving*, North Holland, New York, N.Y., 1979.

Loveland [78]

D. W. Loveland, *Automated Theorem Proving: A Logical Basis*, North-Holland, New York, N.Y., 1978.

Manna and Waldinger [80]

Z. Manna and R. Waldinger, A deductive approach to program synthesis, *ACM Transactions on Programming Languages and Systems*, Vol. 2, No. 1, January 1980, pp. 90-121.

Manna and Waldinger [81]

Z. Manna and R. Waldinger, Deductive synthesis of the unification algorithm, *Science of Computer Programming*, Vol. 1, 1981, pp. 5-48.

Manna and Waldinger [82]

Z. Manna and R. Waldinger, Special relations in program-synthetic deduction, Technical Report, Computer Science Department, Stanford University, Stanford, Calif., and Artificial Intelligence Center, SRI International, Menlo Park, Calif., March 1982.

Manna, Z., and R. Waldinger [85a]

The Logical Basis for Computer Programming, Addison-Wesley, Reading, Mass., Volume 1: Deductive Reasoning (1985), Volume 2: Deductive Techniques (to appear).

Manna, Z., and R. Waldinger [85b]

The origin of the binary-search paradigm, *Ninth International Joint Conference on Artificial Intelligence*, Los Angeles, August 1985.

Morris [69]

J. B. Morris, E-resolution: extension of resolution to include the equality relation, *International Joint Conference on Artificial Intelligence*, Washington, D.C., May 1969, pp. 287-294.

Murray [82]

N. V. Murray, Completely nonclausal theorem proving, *Artificial Intelligence*, Vol. 18, No. 1, 1982, pp. 67-85.

Robinson [65]

J. A. Robinson, A machine-oriented logic based on the resolution principle, *Journal of the ACM*, Vol. 12, No. 1, January 1965, pp. 23-41.

Robinson [79]

J. A. Robinson, *Logic: Form and Function*, North-Holland, New York, N.Y., 1979.

Stickel [82]

M. E. Stickel, A nonclausal connection-graph resolution theorem-proving program. *National Conference on AI*, Pittsburgh, Pa., 1982, pp. 229-233.

Traugott [85]

J. Traugott, Deductive synthesis of sorting algorithms, Technical Report, Computer Science Department, Stanford University, Stanford, Calif. (forthcoming).

Wos and Robinson [69]

L. Wos and G. Robinson, Paramodulation and theorem proving in first order theories with equality, in *Machine Intelligence 4* (B. Meltzer and D. Michie, editors) American Elsevier, New York, N.Y., 1969, pp. 135-150.